

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-110384

(43)Date of publication of application : 23.04.1999

(51)Int.Cl.

G06F 17/24

G06F 17/21

G06F 17/30

(21)Application number : 10-198038

(71)Applicant : HITACHI LTD

(22)Date of filing : 29.06.1998

(72)Inventor : OKAMOTO TAKUYA
MURATA HIDEKO
TAKAHASHI TORU
YAMAZAKI NORIYUKI
AOYAMA YUKI

(30)Priority

Priority number : 09190716
09195408Priority date : 01.07.1997
22.07.1997

Priority country : JP

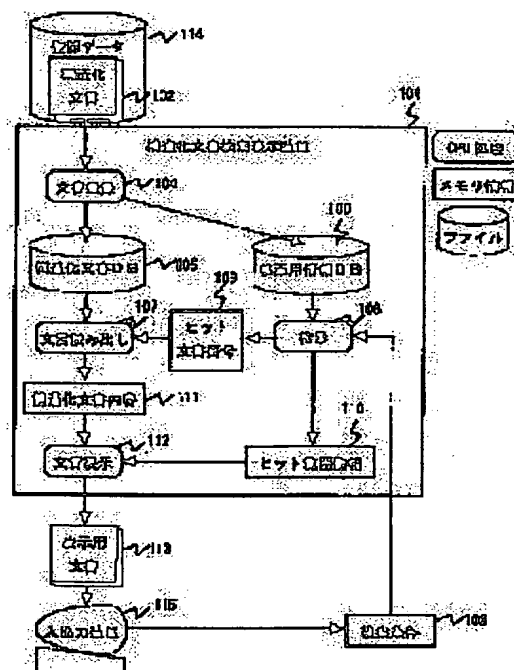
JP

(54) METHOD AND DEVICE FOR RETRIEVING AND DISPLAYING STRUCTURED DOCUMENT

(57)Abstract:

PROBLEM TO BE SOLVED: To retrieve a document that eliminates structure information which becomes an obstacle to retrieval and to display with highlight information added to the original structured document when displaying a retrieval result.

SOLUTION: This device performs document registration processing with a structured document 102 of a file 114 as an input, produces structured document that is undergone a structure analysis and information for document retrieval and stored them in DBs 105 and 106 respectively. Next, when an input-output device 115 inputs a retrieval condition 103, it analyzes the retrieval condition, reads information for document retrieval and performs retrieval processing 108. It outputs document number information 109 that is hit as a retrieval result and hit range information 110. Display processing first reads a corresponding structure-analyzed structured document 111 from the DB 105 based on the information 109 which is hit by document read processing 107. The processing of document display 112 embeds hit information in the document 111 based on the information 110, produces a structured document 113 for display to which the highlight information is added and shows it.



LEGAL STATUS

[Date of request for examination]

21.02.2002

(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開平11-110384

(43)公開日 平成11年(1999) 4月23日

(51)Int.Cl.⁶

G 0 6 F 17/24
17/21
17/30

識別記号

F I

G 0 6 F 15/20

15/40

5 5 4 H

5 5 4 N

5 9 0 E

3 4 0

3 7 0 A

審査請求 未請求 請求項の数14 F D (全 48 頁) 最終頁に続く

(21)出願番号 特願平10-198038

(22)出願日 平成10年(1998) 6月29日

(31)優先権主張番号 特願平9-190716

(32)優先日 平 9 (1997) 7 月 1 日

(33)優先権主張国 日本 (J P)

(31)優先権主張番号 特願平9-195408

(32)優先日 平 9 (1997) 7 月 22 日

(33)優先権主張国 日本 (J P)

(71)出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目 6 番地

(72)発明者 岡本 卓哉

神奈川県横浜市都筑区加賀原二丁目 2 番

株式会社日立製作所システム開発本部内

(72)発明者 村田 英子

神奈川県横浜市都筑区加賀原二丁目 2 番

株式会社日立製作所システム開発本部内

(72)発明者 高橋 亨

神奈川県横浜市都筑区加賀原二丁目 2 番

株式会社日立製作所システム開発本部内

(74)代理人 弁理士 笹岡 茂 (外 1 名)

最終頁に続く

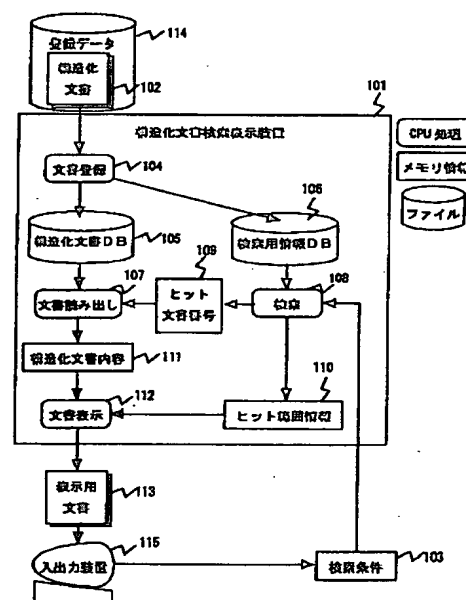
(54)【発明の名称】 構造化文書検索表示方法及び装置

(57)【要約】

【課題】 検索の障害になる構造情報を除去した文書に対して検索し、検索結果の表示の際には元の構造化文書に対してハイライト 情報を付加した表示をする。

【解決手段】 ファイル114の構造化文書102を入力として文書登録の処理を行い、構造解析された構造化文書と、文書検索のための情報を生成し、夫々DB 105、DB 106に格納する。次に入出力装置115から検索条件103が入力されると、検索条件を解析し、文書検索用情報を読み出して検索処理108を行う。検索結果としてヒットした文書番号情報109とヒット 範囲情報110を出力する。表示処理は、まず、文書読み出しの処理107でヒットした文書番号の情報109に基づきDB 105から対応する構造解析済構造化文書111を読み出す。文書表示112の処理では、ヒット 範囲情報110を基に、構造化文書111に対して、ヒット 情報を埋め込み、ハイライト 情報を付加した表示用の構造化文書113を生成し、これを表示する。

【図 1】



<図1は、2の基礎ブロック図>

【 特許請求の範囲】

【 請求項1 】 処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であって、

前記処理装置は、

入力された構造化文書を解析して解析済み構造化文書を生成し、該解析済み構造化文書を前記ファイル装置に格納し、

該解析済み構造化文書から各構造内の内容文字列情報を取得して文書検索用情報を生成し、前記ファイル装置に格納し、

入力された検索条件により該ファイル装置に格納された文書検索用情報を検索し、該検索条件を満たす内容文字列情報があるか否か判定し、該検索条件を満たすとみなされる内容文字列情報を持つ文書の解析済み構造化文書を取得し、かつ該文書の検索条件を満たす範囲の情報を取得し、

該文書の検索条件を満たす範囲をハイライト表示するための表示用文書型定義(表示用DTD)を作成し、

前記文書の検索条件を満たす範囲の情報と表示用文書型定義に基づき構造化文書中にハイライト表示するための情報を付加した表示用構造化文書を作成することを特徴とする構造化文書検索表示方法。

【 請求項2 】 処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であって、

前記処理装置は、

入力された構造化文書を解析して解析済み構造化文書を生成し、該解析済み構造化文書を前記ファイル装置に格納し、

前記入力された構造化文書から予め与えられた検索対象外の構造情報を除去した文書検索用の正規化処理済み構造化文書を生成し、かつ該除去された構造情報を復元するための復元情報を生成し、前記ファイル装置に格納し、

入力された検索条件により該ファイル装置に格納された正規化処理済み構造化文書を検索し、該検索条件を満たす正規化処理済み構造化文書があるか否か判定し、該検索条件を満たすとみなされる文書の正規化処理済み構造化文書を取得し、かつ該文書の検索条件を満たす範囲の情報を取得し、

該文書の検索条件を満たす範囲をハイライト表示するための表示用文書型定義を作成し、

前記検索により取得された正規化処理済み構造化文書を前記復元情報により、除去された構造情報を有する構造化文書に復元し、前記文書の検索条件を満たす範囲の情報と表示用文書型定義に基づき該復元された構造化文書中にハイライト表示するための情報を付加した表示用構造化文書を作成することを特徴とする構造化文書検索表示方法。

【 請求項3 】 処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であって、

前記処理装置は、

入力された構造化文書を解析して解析済み構造化文書を生成し、該解析済み構造化文書を前記ファイル装置に格納し、

該解析済み構造化文書から各構造内の内容文字列情報を取得して文書検索用情報を生成し、前記ファイル装置に格納し、

入力された検索条件により該ファイル装置に格納された文書検索用情報を検索し、該検索条件を満たす内容文字列情報があるか否か判定し、該検索条件を満たすとみなされる内容文字列情報を持つ文書の解析済み構造化文書を取得し、かつ該文書の検索条件を満たす範囲の情報を取得し、

入力された表示対象の部分構造を取得し、

該表示対象の部分構造中の前記検索条件を満たす範囲をハイライト表示するための部分構造表示用文書型定義を作成し、

該表示対象の部分構造に対して、前記文書の検索条件を満たす範囲の情報と部分構造表示用文書型定義に基づき構造化文書中にハイライト表示するための情報を付加した部分構造表示用構造化文書を作成することを特徴とする構造化文書検索表示方法。

【 請求項4 】 請求項1 または請求項2 または請求項3記載の構造化文書検索表示方法において、

検索結果のハイライト表示は、検索ターム毎に複数のハイライト表示形態のいずれかをを用いてハイライト表示することを特徴とする構造化文書検索表示方法。

【 請求項5 】 請求項1 または請求項2 または請求項3記載の構造化文書検索表示方法において、

検索条件中の2つの検索タームについて、各検索タームの相対的な出現位置に関する条件を満たしている場合は、検索条件を構成する各検索タームに対するハイライト表示とその2つの検索タームを含む最小の文字列範囲に対するハイライト表示をそれぞれ異なったハイライト表示形態を用いてハイライト表示することを特徴とする構造化文書検索表示方法。

【 請求項6 】 請求項1 または請求項2 または請求項3記載の構造化文書検索表示方法において、

検索条件に複数の検索タームについて、検索条件を構成する各検索タームに対するハイライト表示と該検索タームを含む構造全体に対するハイライト表示をそれぞれ異なった表示形態を用いてハイライト表示することを特徴とする構造化文書検索表示方法。

【 請求項7 】 請求項4記載の構造化文書検索表示方法において、

各検索ターム毎のハイライト表示のハイライト表示形態は、各検索タームの出現頻度の情報に基づき決定するこ

とを特徴とする構造化文書検索表示方法。

【請求項8】 請求項4記載の構造化文書検索表示方法において、
各検索ターム毎のハイライト表示のハイライト表示形態は、各検索タームごとに予め与えられた重み付けの情報に基づき決定することを特徴とする構造化文書検索表示方法。

【請求項9】 処理装置と、記憶装置と、ファイル装置と、入出力装置を備える構造化文書検索表示装置であって、
前記処理装置は、

入力された構造化文書を解析して解析済み構造化文書を生成し、該解析済み構造化文書を前記ファイル装置に格納する手段と、

前記入力された構造化文書から予め与えられた検索対象外の構造情報を除去した文書検索用の正規化処理済み構造化文書を生成し、前記ファイル装置に格納する手段と、

該除去された構造情報を復元するための復元情報を生成し、前記ファイル装置に格納する手段と、

入力された検索条件により該ファイル装置に格納された正規化処理済み構造化文書を検索し、該検索条件を満たす正規化処理済み構造化文書があるか否かを判定し、該検索条件を満たすとみなされる文書の正規化処理済み構造化文書を取得し、かつ該文書の検索条件を満たす範囲の情報を取得する手段と、

該文書の検索条件を満たす範囲をハイライト表示するための表示用文書型定義を作成する手段と、

前記検索により取得された正規化処理済み構造化文書を前記復元情報により、除去された構造情報を有する構造化文書に復元する手段と、前記文書の検索条件を満たす範囲の情報と表示用文書型定義に基づき該復元された構造化文書中にハイライト表示するための情報を付加した表示用構造化文書を作成する手段を有することを特徴とする構造化文書検索表示装置。

【請求項10】 処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であって、

前記処理装置は、

入力された特定の文書型定義に従う構造化文書をタグを残したままプレーンテキストとして前記ファイル装置に格納し、

入力された検索条件により該ファイル装置に格納されたプレーンテキストを検索し、該検索条件を満たす範囲があるか否かを判定し、該検索条件を満たす範囲を持つ文書をプレーンテキストとして取得し、かつ該文書の検索条件を満たす範囲の情報を取得し、

前記特定の文書型定義を表示用文書型定義とし、前記入力された構造化文書に対して前記検索条件を満たす範囲に対して該表示用文書型定義に基づくハイライト表示す

るための情報を付加した表示用構造化文書を作成することを特徴とする構造化文書検索表示方法。

【請求項11】 請求項10記載の構造化文書検索表示方法において、

検索条件を満たす範囲が構造化文書において文書構造を示すタグの属性情報中に存在するか否かを判定し、

該検索条件を満たす範囲がタグの属性情報中に存在する場合は、構造化文書の内容文字列中に該検索条件を満たす範囲の文字列を含む文字列を追加し、該文字列において該検索条件を満たす範囲に対して前記特定の文書型定義に基づくハイライト表示するための情報を付加した表示用構造化文書を作成することを特徴とする構造化文書検索表示方法。

【請求項12】 請求項10記載の構造化文書検索表示方法において、
入力された検索条件により該ファイル装置にタグを残したままプレーンテキストとして格納された構造化文書を検索する際に、予め指定された特定のタグを構成する文字列を検索対象から除去し、該特定のタグを構成する文字列の前後を連結した文字列に対して検索することで得られる検索条件を満たす範囲に対して、前記特定の文書型定義に基づくハイライト表示するための情報を付加した表示用構造化文書を作成することを特徴とする構造化文書検索表示方法。

【請求項13】 請求項10記載の構造化文書検索表示方法において、

入力された検索条件により該ファイル装置にプレーンテキストとして格納された構造化文書を検索する際に、検索条件を満たす範囲が予め指定された文書構造の開始を示す特定のタグと文書構造の終わりを示す特定のタグに挟まれるか否かを判定し、
挟まれる場合は、文書構造の開始を示す特定のタグより前もしくは文書構造の終わりを示すタグより後ろの内容文字列中に、該検索条件を満たす範囲の文字列を含む文字列を追加し、該文字列において該検索条件を満たす範囲に対して前記特定の文書型定義に基づくハイライト表示するための情報を付加した表示用構造化文書を作成することを特徴とする構造化文書検索表示方法。

【請求項14】 請求項10記載の構造化文書検索表示方法において、

入力された検索条件により該ファイル装置にプレーンテキストとして格納された構造化文書を検索する際に、検索条件を満たす範囲が予め指定された文書構造の開始を示す特定のタグと文書構造の終わりを示す特定のタグに挟まれるか否かを判定し、

挟まれる場合は、文書構造の開始を示す特定のタグより前もしくは文書構造の終わりを示すタグより後ろの内容文字列中に、該検索条件を満たす範囲の文字列を含む文字列を追加し、該文字列において該検索条件を満たす範囲に対して前記特定の文書型定義に基づくハイライト表示するための情報を付加した表示用構造化文書を作成することを特徴とする構造化文書検索表示方法。

【請求項15】 請求項1または請求項2または請求項3または請求項10記載の構造化文書検索表示方法において、

前記表示用文書型定義に基づくハイライト表示するための情報を付加した表示用構造化文書を作成する際に、ハイライト表示するための情報は、検索条件中に指定された方法を用いて付加することを特徴とする構造化文書検索表示方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、SGML、HTMLなどによって作成された構造化文書に対する検索表示

技術に係り、特に構造化文書に対して検索を行い、検索結果に対してハイライトして表示する構造化文書検索表示方法および装置に関する。

【0002】

【従来の技術】ワードプロセッサなどの普及により、作成される文書情報の電子化が進んでいる。これらの電子化文書は、作成される機器、ソフトウェアによって個々のフォーマットを持っており、別の機器あるいはソフトウェアでは、利用できない、あるいは、何らかの変換手段を用意することが必要となっていた。このような文書交換のための共通フォーマットとして、各種の構造化文書が提案されている。これらの構造化文書は、文書の基本構造である、章、節、項などの階層構造を定義できるだけでなく、レイアウト情報を含む事も可能となっている。

【0003】構造化文書の記述言語として、標準化が進められているのが、SGML (Standard Generalized Markup Language) =「標準一般化マークアップ言語」である。SGMLは、構造化文書の構造情報をタグと呼ばれる特定の文字列をテキスト中に埋め込むことで、文書の構造を表現する方法を用いている。SGMLでは、タグの名称、内容、さらに、タグによって示される文書構造をDTD (Document Type Definition) =「文書型定義」によって規定することができる。上記のSGML、DTDについては、「実践SGML」(SGML懇談会実用化WG監訳 1992年4月20日 財団法人日本規格協会発行)に詳細に説明されている。これらの構造化文書を検索システムのDBに登録して、構造名を指定して検索しようとする場合を想定する。登録しようとする各文書のDTDが異なる場合、処理方法としては、文書ごとに文書構造を解析して、指定された構造名がどの部分に相当するかを解析した上で、検索対象の文字列を取得して検索する方法が考えられる。しかし、この方法は、多くの処理時間を必要とする。また、構造名ごとに各文書の対応する箇所をテーブルで持つなどの方法を用いる場合、各文書に出現する構造名を全て一括して管理し、構造名ごとに各文書の対応する部分を登録する必要がある、膨大な管理テーブルが必要となる。さらに、異なるDTDが混在する文書を登録しても、検索対象の構造をすべての文書が持っているとは限らず、また、例えば、「要約」、「要旨」のように、同じ内容であっても異なる構造名を付けた場合、これらの異なる構造名を全て指定して、検索を行なわなければならない、現実的な構造化文書の検索とは考えられない。

【0004】したがって、同じ文書型定義で生成された文書だけを登録するように運用することが構造化文書の検索では必要となる。あらかじめ指定された構造名について、各文書の対応する部分を管理する。検索の際には、検索対象の構造名および検索条件を指定すると、各文書の指定された構造に対応する部分に検索条件に当て

はまる文字列が含まれると、検索条件にヒットしたと判断される。

【0005】構造化文書の検索結果として、文書の内容を表示するための機能の従来技術について以下に述べる。まず、第1の従来技術として、特開平8-339369「文書表示装置および文書表示方法」が挙げられる。本従来技術は、SGML文書の構造解析および構造表示用のレイアウトへの変換、さらに指定構造の内容の表示を行う方法について述べられており、本技術を用いることで構造化文書を構造単位で表示することが可能である。さらに、本従来技術においては、指定構造のハイライト表示(強調した表示のことであり、色、字体、字の大きさ等を変えたり、アンダーラインを付したりする)の手段を提供している。しかし、ここで示されているハイライト表示手段とは、構造毎に表示方法をコントロールする手段であり、構造単位に、表示の有無、ハイライト表示などの指定を行う。したがって、本従来技術において、構造化文書の検索結果の表示を実現する際に必要となる、ヒットした検索タームに対するハイライト表示を実現する方法が示されているわけではない。

【0006】また、第2の従来技術としては、特開平8-212230「文書検索方法および文書検索装置」で構造化文書以外の文書の検索結果に対するハイライト表示方法が示されている。しかし、本従来技術は、表示するためのテキストに対するヒット範囲の取得およびハイライト情報の付加を実現するのみであり、構造化文書の検索結果として得られた文書に対してハイライト情報を付加する機能を持つわけではない。

【0007】上記2つの従来技術を組み合わせただけでは、構造化文書に対する検索結果として出力する文書に対して、ヒットしたタームに対するハイライト情報の付加を実現する事はできない。つまり、構造化文書において、ハイライト表示を実現するためには、表示対象の文書の作成時のDTDにハイライト用の構造情報を追加したDTDを作成する手段が必要となる。

【0008】構造化文書にハイライト情報を付加した際の文書型定義の変更方法については、第3の従来技術である、特願平8-159202「構造化文書の版管理方法および装置」に、元のDTDに対して新たな構造を追加したDTDを生成する方法が示されている。本従来技術を用いる事により、ハイライト情報を付加した文書型定義を作成することができる。

【0009】第1、第2の従来技術により、構造化文書を構造が分かるように表示すること、さらに構造化されていない文書においては、ヒット範囲のハイライト表示をする事が可能であることがわかる。さらに、第3の従来技術を用いることにより、構造ごとに取得したハイライト情報を付加した文書型定義が指定できる。これらの技術を組み合わせることで、構造化文書の特定の構造の検索結果に対してハイライト情報を付加した構造化文書

を出力し、ハイライト表示を実現する事ができる。

【0010】また、最新の情報を入手する方法として、近年インターネットが爆発的に広まっている。インターネット上に存在する数多くの情報から自分が必要とする情報をいち早く知る手段として、Web上の情報の検索機能も充実してきた。HTML (Hyper Text Markup Language) は、WWW (World Wide Web) 上において、文書内容を記述し、他の資源へのリンク情報、文書のフォーマットを表現するための言語である。HTMLは、特定のDTDにしたがって記述されたSGMLとみなすことができる。このHTML文書を作成、加工する手段として、HTMLエディタがある。また、作成されたHTML文書を解析し、表示するHTMLブラウザが存在する。HTMLブラウザには、検索する文字列(以下、「検索ターム」という。)を入力し、表示中のHTML文書に対して検索を行い、ヒットした箇所を反転表示などの強調表示を行う機能を持つものがある。SGMLについても、レイアウト表示し、加工する機能を持つSGMLブラウザが存在する。SGMLブラウザには、ブラウザ上に表示中のSGML文書に対して、全文検索し、検索条件に適合する箇所をハイライト表示する。これらのブラウザでは、文書表示の際に文書の解析を行ない、表示用のデータを作成している。検索はこのブラウザ上の表示用のデータに対して検索を行ない、画面上でヒット位置をハイライト表示している。

【0011】

【発明が解決しようとする課題】上記の従来技術の組み合わせにより、与えられた構造化文書に対して、構造毎に検索した結果を、個々にハイライト表示することが可能である。しかし、構造情報には、章、節、項のように文書構造そのものを表わしているものだけでなく、アンダーラインの付加などレイアウト用の情報も含まれる場合がある。これらの構造情報は、必ずしも文の切れ目で挿入されるとは限らない。文書検索する際には、このような構造情報を除去しなければ、文書中に含まれている語であるにもかかわらず、検索できないという問題がある。このように、検索時に不要となる構造情報を除去する処理を、以下の説明では「正規化処理」と呼ぶ。正規化処理を行なった構造化文書を検索対象とし、元の構造化文書に対してハイライト情報を付加した表示を実現するためには、正規化処理を行なった構造化文書に対して、上記の従来技術を用いた方法を利用するだけでは実現できない。つまり、この方法では、検索時には、元の文書の構造情報の一部しか残っていないため、この構造情報に対してハイライト情報を付加するだけでは、元の構造化文書に対してヒットした検索タームのハイライト表示を実現することにならないのである。

【0012】一方、HTML文書は、ブラウザ依存の独自の拡張により複数のDTDに基づいて作成されたHTML文書が存在し、またどのDTDに基づいて記述され

ているかがわからない。さらに、SGMLの文法に基づいて正しく記述されていない文書も数多く存在するため、SGMLと同様の方法で構造解析することは困難である。また、(1)プレーンなテキスト文書に対しては、検索処理を行い、検索ヒット位置の前後にハイライト用のタグを挿入したHTML文書を生成することにより、HTMLブラウザ上で、検索ヒットした文字列を強調表示することが可能である。しかし、タグ内の文字列が検索タームと一致した場合、この検索ヒット位置の前後に対して、ハイライト用のタグを挿入すると、元々のHTMLのタグの内容が変更されるため、正しく表示されなくなるといった問題が起こる。さらに、(2)HTMLブラウザ上で連続して表示されている文字列の途中に、レイアウトを表現するタグが挿入されている場合があり、HTML文書に対して検索する場合は、タグを除いて検索しなければ正しく検索することができない。例えば、HTML文書中に「今月の特集記事」と書かれており、検索タームを「特集記事」とした場合、HTML文書中では、「特集」と「記事」の間に文字を拡大して表示するための「」のタグが記述されているため、タグを飛ばして検索しなければ正しく検索することができない。

【0013】本発明の目的は、正規化処理された文書に対する検索結果から、元の文書に対するハイライト情報の付加を実現するために、検索用の文書から、元の文書のハイライト範囲情報への変換を実現することにある。本発明の他の目的は、正規化後のヒットタームが、元の文書において複数の構造にまたがっている場合、各構造ごとに、ヒットした範囲に対してハイライト情報を付加し、ハイライト表示することにある。本発明のさらに他の目的は、ヒットしたタームが含まれる構造全体のハイライト表示、あるいは、出現位置の距離条件を満たした2つの検索タームを含む領域全体をハイライト表示するなどの処理をするため、階層的なハイライト情報を付加し、異なるハイライト表示形態によりハイライト表示することにある。本発明のさらに他の目的は、構造化文書の部分構造だけを抽出して表示する場合に、このような部分構造の内容についても、ハイライト情報を付加し、ハイライト表示することにある。本発明のさらに他の目的は、文書構造を示すHTMLタグが存在する文書から文字列を検索する場合、設定した検索タームと一致した文字列がHTMLタグ内に存在する場合や、検索タームがHTMLタグをまたがって記述されている場合でも検索を可能にすることにある。本発明のさらに他の目的は、検索条件にヒットした文字列をハイライト表示可能にすることにある。

【0014】

【課題を解決する為の手段】上記の課題を解決するため、本発明は、処理装置と、記憶装置と、ファイル装置

と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であり、前記処理装置は、入力された構造化文書を解析して解析済み構造化文書を生成し、該解析済み構造化文書を前記ファイル装置に格納し、該解析済み構造化文書から各構造内の内容文字列情報を取得して文書検索用情報を生成し、前記ファイル装置に格納し、入力された検索条件により該ファイル装置に格納された文書検索用情報を検索し、該検索条件を満たす内容文字列情報があるか否か判定し、該検索条件を満たすとみなされる内容文字列情報を持つ文書の解析済み構造化文書を取得し、かつ該文書の検索条件を満たす範囲の情報を取得し、該文書の検索条件を満たす範囲をハイライト表示するための表示用文書型定義(表示用DTD)を作成し、前記文書の検索条件を満たす範囲の情報と表示用文書型定義に基づき構造化文書中にハイライト表示するための情報を付加した表示用構造化文書を作成するようにしている。

【0015】また、処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であり、前記処理装置は、入力された構造化文書を解析して解析済み構造化文書を生成し、該解析済み構造化文書を前記ファイル装置に格納し、前記入力された構造化文書から予め与えられた検索対象外の構造情報を除去した文書検索用の正規化処理済み構造化文書を生成し、かつ該除去された構造情報を復元するための復元情報を生成し、前記ファイル装置に格納し、入力された検索条件により該ファイル装置に格納された正規化処理済み構造化文書を検索し、該検索条件を満たす正規化処理済み構造化文書があるか否か判定し、該検索条件を満たすとみなされる文書の正規化処理済み構造化文書を取得し、かつ該文書の検索条件を満たす範囲の情報を取得し、該文書の検索条件を満たす範囲をハイライト表示するための表示用文書型定義を作成し、前記検索により取得された正規化処理済み構造化文書を前記復元情報により、除去された構造情報を有する構造化文書に復元し、前記文書の検索条件を満たす範囲の情報と表示用文書型定義に基づき該復元された構造化文書中にハイライト表示するための情報を付加した表示用構造化文書を作成するようにしている。

【0016】また、処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であり、前記処理装置は、入力された構造化文書を解析して解析済み構造化文書を生成し、該解析済み構造化文書を前記ファイル装置に格納し、該解析済み構造化文書から各構造内の内容文字列情報を取得して文書検索用情報を生成し、前記ファイル装置に格納し、入力された検索条件により該ファイル装置に格納された文書検索用情報を検索し、該検索条件を満たす内容文字列情報があるか否か判定し、該検索条件を満たすとみなされる内容文字列情報を持つ文書の解析済

み構造化文書を取得し、かつ該文書の検索条件を満たす範囲の情報を取得し、入力された表示対象の部分構造を取得し、該表示対象の部分構造中の前記検索条件を満たす範囲をハイライト表示するための部分構造表示用文書型定義を作成し、該表示対象の部分構造に対して、前記文書の検索条件を満たす範囲の情報と部分構造表示用文書型定義に基づき構造化文書中にハイライト表示するための情報を付加した部分構造表示用構造化文書を作成するようにしている。

【0017】また、処理装置と、記憶装置と、ファイル装置と、入出力装置を備える構造化文書検索表示装置であり、前記処理装置は、入力された構造化文書を解析して解析済み構造化文書を生成し、該解析済み構造化文書を前記ファイル装置に格納する手段と、前記入力された構造化文書から予め与えられた検索対象外の構造情報を除去した文書検索用の正規化処理済み構造化文書を生成し、前記ファイル装置に格納する手段と、該除去された構造情報を復元するための復元情報を生成し、前記ファイル装置に格納する手段と、入力された検索条件により該ファイル装置に格納された正規化処理済み構造化文書を検索し、該検索条件を満たす正規化処理済み構造化文書があるか否か判定し、該検索条件を満たすとみなされる正規化処理済み構造化文書の情報を取得し、かつ該文書の検索条件を満たす範囲の情報を取得する手段と、該文書の検索条件を満たす範囲をハイライト表示するための表示用文書型定義を作成する手段と、前記検索により取得された正規化処理済み構造化文書を前記復元情報により、除去された構造情報を有する構造化文書に復元する手段と、前記文書の検索条件を満たす範囲の情報と表示用文書型定義に基づき該復元された構造化文書中にハイライト表示するための情報を付加した表示用構造化文書を作成する手段を有するようにしている。

【0018】また、処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であり、前記処理装置は、入力された特定の文書型定義に従う構造化文書をタグを残したままプレーンテキストとして前記ファイル装置に格納し、入力された検索条件により該ファイル装置に格納されたプレーンテキストを検索し、該検索条件を満たす範囲があるか否か判定し、該検索条件を満たす範囲を持つ文書をプレーンテキストとして取得し、かつ該文書の検索条件を満たす範囲の情報を取得し、前記特定の文書型定義を表示用文書型定義とし、前記入力された構造化文書に対して前記検索条件を満たす範囲に対して該表示用文書型定義に基づくハイライト表示するための情報を付加した表示用構造化文書を作成するようにしている。

【0019】また、処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であり、前記処理装置は、入力された特定の文書型定義に従う構造化文書をタグを残し

たままプレーンテキストとして前記ファイル装置に格納し、入力された検索条件により該ファイル装置に格納されたプレーンテキストを検索し、該検索条件を満たす範囲があるか否かを判定し、該検索条件を満たす範囲を持つ文書をプレーンテキストとして取得し、かつ該文書の検索条件を満たす範囲の情報を取得し、検索条件を満たす範囲が構造化文書において文書構造を示すタグの属性情報中に存在するか否かを判定し、該検索条件を満たす範囲がタグの属性情報中に存在する場合は、構造化文書の内容文字列中に該検索条件を満たす範囲の文字列を含む文字列を追加し、該文字列において該検索条件を満たす範囲に対して前記特定の文書型定義に基づくハイライト表示するための情報を付加した表示用構造化文書を作成するようにしている。

【0020】また、処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であり、前記処理装置は、入力された特定の文書型定義に従う構造化文書をタグを残したままプレーンテキストとして前記ファイル装置に格納し、予め指定された特定のタグを構成する文字列を検索対象から除去し、該特定のタグを構成する文字列の前後を連結した文字列に対して検索することで得られる検索条件を満たす範囲に対して、前記特定の文書型定義に基づくハイライト表示するための情報を付加した表示用構造化文書を作成するようにしている。

【0021】また、処理装置と、記憶装置と、ファイル装置と、入出力装置を備える情報処理システムにおける構造化文書検索表示方法であり、前記処理装置は、入力された特定の文書型定義に従う構造化文書をタグを残したままプレーンテキストとして前記ファイル装置に格納し、入力された検索条件により該ファイル装置にプレーンテキストとして格納された構造化文書を検索する際に、検索条件を満たす範囲が予め指定された文書構造の開始を示す特定のタグと文書構造の終わりを示す特定のタグに挟まれるか否かを判定し、挟まれる場合は、文書構造の開始を示す特定のタグより前もしくは文書構造の終わりを示すタグより後ろの内容文字列中に、該検索条件を満たす範囲の文字列を含む文字列を追加し、該文字列において該検索条件を満たす範囲に対して前記特定の文書型定義に基づくハイライト表示するための情報を付加した表示用構造化文書を作成するようにしている。

【0022】

【発明の実施の形態】第1の実施例の概略の処理ブロック図を図1に示す。101は、構造化文書検索表示装置である。登録データファイル(114)に格納された、構造化文書(102)を入力として文書登録の処理を行う事で、構造解析された構造化文書(図3により後述する)と、文書検索のための文書検索用情報(図5により後述)が生成される。構造解析された構造化文書は、構造化文書データベース(以下、データベースをDBと記

述する。)(105)に格納し、検索用情報は、検索用情報DB(106)に格納される。次に入出力装置(115)から、検索条件(103)が入力されると、検索条件を解析し、文書検索用情報を読み出して、検索処理(108)を行う。検索結果としては、ヒットした文書番号の情報(109)とヒット範囲の情報(110)を出力する。表示処理は、まず、文書読み出しの処理(107)で、ヒットした文書番号の情報(109)に基づいて、構造化文書DB(105)から、指定された構造解析済構造化文書(111)を読み出す。文書表示(112)の処理では、ヒット範囲情報(110)を基に、構造解析済構造化文書(111)に対して、ヒット情報を埋め込んだ表示用の構造化文書(113)を生成する。生成された表示用の構造化文書は、入出力装置(115)に表示される。

【0023】図2に構造化文書検索表示の処理フローを示す。まず、構造化文書の登録処理を行なう(201)。登録処理の内容については、図4のフローチャートを用いて後述する。次に、指定された検索条件を用いて構造化文書を検索する(202)。検索処理の詳細は、図6のフローチャートを用いて後述する。検索結果としては、ヒット文書数とヒット文書を識別する番号と各文書毎の検索タームのヒット範囲がある。ヒット範囲の情報は、ヒットした検索タームが含まれる構造を識別するための構造ID(構造識別子)と構造内でのヒット開始位置、テキスト長の情報を出力する。構造化文書検索の処理で、ヒット文書数が1以上であれば(203)、順次、ヒットした文書の内容を読み出し(204)、読み出した文書のヒット範囲情報を取得し(205)、ハイライト表示を実現する(206)。表示処理の詳細については、図9を用いて後述する。さらにヒットした文書があれば、204から206の処理を繰り返す。表示処理を終えると、次の検索処理の有無を確認し(208)、検索条件がなければ、処理を終え、検索条件があれば、202の処理に戻って構造化文書の検索表示処理を繰り返す。

【0024】図3は、構造化文書登録処理の概要を示した図である。まず、SGML文書(301)の構造を解析し、木構造(302)を生成する。生成した木構造の各項目の内容をテーブル形式のデータ(303)として出力し、これを解析済み構造化文書として登録する。ここで、CDATAとは文字列データのことである。

【0025】図4は構造化文書登録処理のフローチャートである。まず、構造化文書を解析する(401)。解析された構造化文書を解析済構造化文書として登録する(402)。構造化文書の解析には、DTDを利用してSGML文書を解析するSGMLパーサを用いることで実現できる。次に、解析された構造化文書に対して、検索に不要な構造を除去するための正規化処理を行なう(403)。正規化処理の手順については、図12を用

いて後述する。そして、正規化処理した構造化文書を、文書データベースに登録する(404)。さらに、データベースに登録された解析済み構造化文書から、構造化文書の検索に必要な検索用情報として、構造情報、構造内のテキストの情報を取り出す(405)。ここで得られた検索用情報を検索用情報DB(106)に登録する(406)。ここで、登録される検索用情報は、SGML文書中の構造情報(タグ)を除去し、各構造ごとに構造情報とその内容を表すテキスト列を格納したものである。図5に上記検索用情報と正規化した構造化文書からなる検索用のテキストの格納例を示す。上記処理を登録文書に対して繰り返し実行し、登録文書が無くなったとき処理を終了する(407)。登録内容は登録文書の全文検索に用いる。図5は、検索用のテキストとして、出力される内容の例である。このように文書構造の構造IDとテキスト列を対応付けるテーブルと文字列情報からなる情報を検索用のテキストとして登録する。検索の際には、構造IDを元に必要な文字列を抽出して検索を行なう。

【0026】図6は、図2の構造化文書検索表示処理の202ステップの構造化文書検索の処理フローである。検索条件は、「検索対象の構造指定: 検索条件式」のように与えられる。検索対象の構造は、例えば、「<文書.タイトル>」のように、「<」と「>」で囲まれ、上位構造(例の場合、「文書」と下位構造(例の場合、「タイトル」)は「.」で区切られ、階層構造中のどの構造に対して、検索を行なうかが指定される。検索条件式は、例えば、and("検索","文書")では、「検索」と「文書」が両方出現する条件を示しており、C<=10("検索","文書")では、「検索」と「文書」が10文字以下の文字を挟んで出現する条件を示している。

【0027】構造化文書検索は、まず、ヒット文書数のカウンタをクリアし(601)、次に、検索条件中の検索対象の構造指定の部分の解析を行なう(602)。ステップ602では、<文書.タイトル>のように構造を指定する文字列から、解析済み構造化文書の対応する構造を一意に特定できる構造ID(構造識別子)を取得する。構造ID取得の処理内容は、図7のフローチャートを用いて後述する。次に、検索対象として登録された文書(検索用のテキスト)を読み出し、ステップ602で取得した指定構造IDに対応するテキスト部分を取得する(603)。検索条件から、検索ターム、さらに複数の検索タームの出現の論理積、距離条件などの論理条件からなる検索条件式を解析し(604)、得られた検索タームによりステップ603で取得されたテキスト部分の全文検索を行ない、検索条件式の論理条件を満たすか否かの判定、すなわち、検索条件にヒットしたか否かを判定する(605)。検索条件にヒットすると(606)、検索結果として文書の番号、検索タームが含まれる構造のIDと、構造中の検索タームがヒットした範囲

の情報を出力する(607)。さらに、ヒットした文書の数のカウントし(608)、本処理を全文書について行なった後(609)、ヒット文書数を出力する(610)。

【0028】図7は、図6の検索条件の解析における、構造指定内容の解析処理のフローチャートである。まず、文書の最上位構造を取得する(701)。次に最上位構造から順に下位構造を取得する。取得した構造が指定構造の下位構造であれば(703)、その構造を検索対象の構造として構造IDを出力する(704)。下位構造があれば(705)、さらにその下位構造に対して、同様に指定された構造の下位構造か否かを判定し、下位構造であれば構造IDを出力する処理(706)を下位構造がなくなるまで繰り返し(707)、全ての構造について処理が終われば、検索対象の構造IDの一覧が得られる。図8に検索対象となる構造ID一覧の出力形式を示す。検索対象となる構造IDの数(801)と、検索対象として得られた数のID(802)が出力される。

【0029】図9は、表示処理の内容を示すフローチャートである。本フローチャートを用いて、表示処理の内容を以下に述べる。まず、検索対象の構造化文書は、検索に不要な構造を除去する正規化処理を行なった後の文書であるため、検索によりヒットした構造およびヒット範囲情報は、必ずしも登録した正規化していない構造化文書における構造および範囲と一致するとは限らない(図3の木構造302と図12の木構造1203を参照)。表示に用いる文書は、登録した正規化していない構造化文書に対して、ヒットした範囲にハイライト情報を付加した文書となる。したがって、まず、登録文書のDTDから、表示に用いる文書用の表示用DTDの作成処理を行なう(901)。表示用DTD作成処理の内容については、図11を用いて後述する。さらに、正規化後の構造化文書に対して得られたヒット範囲については、正規化前の登録した構造化文書における構造およびハイライト範囲情報に変換する(902)。正規化後の文書のヒット範囲情報の正規化前の文書のハイライト範囲情報への変換処理の内容については、図15を用いて後述する。

【0030】次に表示に用いる解析済み文書の最上位構造の情報を読み出し、903から911の処理を順に繰り返すことで、表示用の文書の出力処理を行なう。まず、構造情報を読み出し(903)、最初に構造の開始タグを出力する(904)。さらに本構造に下位構造が存在するなら(905)、下位構造に対して、表示処理(903から911の処理)を再帰的に行なう(906)。下位構造がなくなれば、構造の終わりを示すタグを出力する処理(911)に移る。

【0031】ここで、下位構造とは、文字列を含む。したがって、

<文書>

<タイトル>

構造化文書

</タイトル>

<本文>

<強調>構造化文書</強調>の検索は、・・・

</本文>

</文書>

などの構造化文書については、<タイトル>の下位構造として、文字列(SGMLでは、CDATAと表現される)という構造が存在することになる。CDATAは、下位構造を持たず、文字列情報として、上記の例の場合、「構造化文書」という内容を持つのである。<本文>についても同様に、<強調>という構造と、「の検索は、・・・」という内容を持つ文字列が下位構造として存在することになる。

【0032】905のステップで下位構造が存在しないと判定された場合は、文字列の構造であるため、本構造の内容に対して、ヒット範囲情報と比較し(908)、ヒット範囲が含まれる構造であれば、ハイライト処理を行なう(909)。ハイライト処理については、図16を用いて後述する。ヒット範囲が含まれない文字列であれば、内容をそのままテキストとして出力する(910)。出力内容が文字列の場合は、904、911のステップで、開始タグ、終了タグは出力しない。上記の処理で構造ごとのハイライト表示を実現する。さらに処理すべき構造があれば、903からの処理を繰り返す(912)。

【0033】図10は、登録用DTD(1001)と、登録するSGML文書(文書インスタンス)の例(1002)、ハイライト表示に用いる表示用DTD(1003)と、表示用に変換したSGML文書(文書インスタンス)の例(1004)である。なお、DTD(Document Type Definition)とは、従来の技術の項で述べたように、タグの名称、内容、さらに、タグによって示される文書構造を規定する文書型定義である。DTDにおいて、構造を表現する場合は、"<!ELEMENT タグ名"に続いて、"-または"O"が2つ並べられる。最初の"-または、"O"は、構造開始タグの省略の可否を示しており、"-の場合は、省略できない。"O"の場合は省略可能である。2つめの"-または"O"は、終了タグの省略の可否を示している。次に、内容モデルとして、下位構造に出現しうる構造が記述される。図10のDTD1001の(タイトル、本文)の場合、タイトルは下位構造1、本文は下位構造2である。"(下位構造1,下位構造2?)"のように記述される場合は、下位構造1の後に下位構造2がそれぞれ1回だけ出現することを示し、"?は、下位構造2は、出現しなくても良いことを示している。"(下位構造1|下位構造2)*"の場合は、下位構造1、2が順序不同で複数回(0回を含む)出現することを示す。こ

で、内容モデルに"CDATA"と記述されている場合は、その構造中には、1つだけの文字列が存在することを示している。#PCDATAも文字列を表わしているが、繰り返し出現が可能である。文字列と、構造が混在する場合は、#PCDATAを用いる必要がある。

【0034】内容モデルに、"CDATA"の代わりに"RCDATA"が指定される場合がある。CDATAとRCDATAの違いは、CDATAが、構造内にエンティティ参照("<xxxx;"のように記述される。外字への置き換えなどに利用される。)が出現した場合に、エンティティ(外字など)への変換を行なわないで、出現した文字列のまま、文字列として扱うのである。"RCDATA"が指定された場合は、エンティティへの変換を行なった文字列を、文字列として扱う。

【0035】ハイライト表示するためには、文字列に対してハイライト情報を付加できるように、文書構造を変更する必要がある。1003にアンダーラインで示した変更点のように、各構造の文字列部分に対しては、全てハイライト表示用の構造情報を追加し、さらにハイライト表示用の構造情報("<!ELEMENT ハイライト --(#PCDATA)"を付加する必要がある。元のDTDで内容モデルの"CDATA"となっている部分が、"(#PCDATA|ハイライト)*"に変更されているのは、CDATAがその構造中には、文字列が1つしか存在しないことを示しており、繰り返しの要素としては出現し得ないためである。ハイライト用のタグが付加されるため、元の構造がCDATAであっても、#PCDATAに変更した上で、ハイライトが繰り返し出現することが可能なように、"(#PCDATA|ハイライト)*"とするのである。

【0036】図11は、登録用のDTDからハイライト表示用のDTDを作成するための処理内容を表すフローチャートである。まず、登録用DTDを読み出し(1101)、DTDの内容を解析して、ELEMENT項目を取得する(1102)。ELEMENT項目の内容モデル中に、CDATA、RCDATA、#PCDATAなどが指定されている場合は、全て、ハイライト用の構造を付加できるように、内容モデルを変更する(1103-1106)。内容モデルの変更は、まず、"CDATA"、"RCDATA"を"#PCDATA"に変更した上で、"#PCDATA"を"(#PCDATA|ハイライト)*"のように、ハイライトタグで囲まれた文字列と、囲まれていない文字列が繰り返し出現するように定義する。元の内容モデルが、"(#PCDATA|アンダーライン)*"のように複数の構造が、繰り返し出現するように記述されている場合は、"(#PCDATA|アンダーライン|ハイライト)*"のように、ハイライト構造が出現することを記述するだけで良い。すべてのELEMENT宣言について変更処理が終わると(1107)、ハイライト用の構造の定義として、"<!ELEMENT ハイライト -- CDATA"を追加する(1108)。以上の処理で、図10の1003に示したハイライト表示用のDTDが生成される。

【0037】図12は、構造化文書の正規化処理の内容

を示した図である。図10の1001に示した構造化文書を木構造に表わすと1201のようになる。不要な構造として"アンダーライン"が指定されている場合、正規化処理の最初の処理として、1202に示すように、アンダーラインという構造を削除し、アンダーラインの下位構造に含まれる文字列は、直接上位構造である"本文"の要素とする。さらに、"本文"の下位構造として、文字列(CDATA)が2つ並んでいるため、1203のように、文字列を連結して、1つの文字列データとする。

【0038】図13は、正規化処理前の構造化文書(1301)、正規化処理後の構造化文書(1302)の内容を解析し、テーブル形式に変換して出力した内容である。1303は、構造情報を格納したテーブルであり、0から6までの構造IDが付けられた構造は、正規化前の構造の情報である。0が最上位構造であり、下位構造の情報をたどっていくことで、文書構造が分かる。7から9までの構造ID(構造識別子)が付けられた構造は、正規化後に変更、追加された構造である。7が最上位構造であり、下位構造を辿ると正規化後の文書構造が分かる。ここで、変更のない構造である"タイトル"以下の構造である構造ID1, 2の構造情報はそのまま残される。さらに、正規化処理で追加された構造ID7から9の構造については、1304の正規化対応テーブルにより、正規化前の構造との対応関係が格納される。

【0039】図14は、正規化後の構造化文書に対して、検索した際のヒット範囲の情報を正規化前の構造化文書における範囲情報へ変換した結果を示している。1401の正規化後の構造情報に基づいて得られたヒット範囲の情報を、図13の1304の正規化対応テーブルの情報を利用して、正規化前の構造化文書における範囲情報(1402)に変換している。本図の例では、正規化後の構造ID9のヒット範囲が、正規化前の文書では、構造ID5と6に分かれているため、2つの構造中のハイライト対象の範囲情報に変更している。

【0040】図15に、図9の902ステップの正規化処理後の構造化文書に対するヒット範囲情報を正規化処理前の構造化文書に対するヒット範囲情報に変換する処理内容のフローチャートを示す。まず、正規化後のヒット範囲情報を順次読み出し(1501)、ヒット範囲情報の構造IDが、正規化後に追加されたものか、正規化前から存在するものであるかを判定する(1502)。正規化前から存在する構造IDであれば変更はないため、そのまま、正規化前のヒット範囲情報として出力する(1503)。正規化後に作成された構造IDであれば、図14の正規化対応テーブルの正規化後構造IDを辿り、文字範囲の情報から、対応する正規化前の構造IDと、ヒット範囲を得る(1504)。正規化処理前の構造におけるヒット範囲を得たら、これを正規化前のヒット範囲として出力する(1505)。全てのヒット範囲情報について処理を終える(1506)と、表示用の

ハイライト範囲情報が得られる。

【0041】図16は、図9の909ステップのハイライト処理のフローチャートである。まず、文書の先頭から、ハイライト開始までの文字列を出力する(1601)。次に、ハイライト表示に用いる構造の開始タグを出力する(1602)。さらに、ハイライト範囲の文字列を出力し(1603)、ハイライト表示に用いる構造の終了タグを出力する(1604)。すべてのハイライト処理を終えると(1605)、残ったテキストを出力し、ハイライト処理を終わる(1606)。

【0042】次に第2の実施例として、ヒット条件によって、ハイライト表示方法を変更する処理、さらに複数のハイライト処理を階層的に行なう場合の処理について説明する。概略処理ブロック図は、図1と同じである。図17は、本実施例で用いるヒット範囲情報(1701)である。図14に示したヒット範囲情報に対して追加された情報は、各ヒットした条件を格納する領域(1702)が追加されていることである。さらに、図14では、ヒットした検索タームの範囲だけを出力しているが、検索条件によって、ヒットした検索タームに加えて、その検索タームが含まれる構造全体に対するハイライトなど、検索タームを含む領域を指定することを可能としている。これらのヒット条件の情報は、構造化文書の検索処理時に付加する。ここでは、検索条件に用いられた距離条件、各検索タームの出現頻度などの情報を付加しているが、検索ターム毎にあらかじめ、重み付けを行なうなどの方法を用いることもできる。

【0043】図18は、ヒット条件とハイライト方法(ハイライト表示形態)の対応を定義したテーブル(1801)である。ヒット条件(1802)に対応するハイライト方法(1803)が記述されている。各ヒット条件によって、ヒットした範囲は、本テーブルの内容に基づいてハイライト表示を行なう。さらに、階層情報(1804)が与えられており、階層情報の値が大きいほど、構造全体のハイライトなど上位のハイライト構造となっている。

【0044】図19は、上記のハイライト処理を実現するための、表示用DTD作成の処理内容を示したものである。登録に用いた元のDTD(1901)に対して、上位のハイライト構造内には下位のハイライト構造を階層的に指定でき、さらに省略も可能なように定義を変更、追加したハイライト表示用のDTD(1902)を生成している。DTDの作成方法は、図11を用いて前述した処理に対して、1106ステップのハイライト情報付加の際に、複数存在するハイライト情報をすべて付加(1903)し、さらに1108ステップのハイライト用ELEMENT宣言追加の際に、図18の階層情報(1804)を元に、各ハイライト構造の下位構造として、下位のハイライト構造および文字列を内容モデルとして持つようにすれば良い。下位のハイライト構造がなけれ

ば、内容モデルとして、文字列だけが出現する(1904)。

【0045】図20は、第2の実施例におけるハイライト処理のフローチャートである。まず、ハイライト情報を開始位置順を第1キー、階層情報の上位から下位の順を第2キーとしてソートする(2001)。次に、ハイライト開始までのテキストを出力し(2002)、ハイライト開始タグを出力する(2003)。さらに、ハイライト範囲の終わりまでに、次のハイライトが開始していれば、下位の構造情報が存在するため(2004)、その位置までのテキストを出力した上で(2005)、下位のハイライト構造におけるハイライト処理を行なう(2006)。下位構造におけるハイライト処理は、2003から2009の処理と同じである。下位のハイライト構造に対する処理を終えた後、さらに下位のハイライト構造があれば(2007)、2005ステップの処理に戻って、次のハイライト構造までのテキストを出力し、下位のハイライト構造の処理を行なう。下位のハイライト構造がなくなれば、構造の終わりまでのテキストを出力して(2008)、ハイライト終了タグを出力する(2009)。ハイライトの情報が残っていれば、2002のステップに戻り、処理を繰り返す。ハイライトの情報が終われば(2010)、残ったテキストを出力し、処理を終える(2011)。

【0046】図21は、上記処理により生成されるSGML文書の例である。図22は、図21のSGML文書の本文の表示例である。重なったハイライト範囲については、複数のハイライトのための表示方法を重複して行なっている。

【0047】第3の実施例として、構造化文書の部分構造だけを切り出し、ハイライト表示する場合の処理内容を示す。図23は、本実施例の概略処理ブロック図を示したものである。図1からの変更点は、表示対象の構造(2301)を指定するようにしていることと、表示対象の構造の指定内容を元に、文書表示(112)の処理の代わりに部分構造表示の処理(2302)を行なっていることである。

【0048】図24は、部分構造を抽出して、表示する場合の処理手順を示したフローチャートである。まず、部分構造表示用のDTDを作成する(2401)。部分構造表示用のDTDの作成処理については、図26を用いて後述する。さらに、正規化後の構造化文書に対して得られたヒット範囲については、正規化前の登録時の文書における、構造IDおよびヒット範囲情報に変換する(2402)。正規化後の文書の情報の正規化前の文書の範囲情報への変換処理の内容については、図16を用いて前述した方法を用いることができる。次に表示対象となっている解析済み文書の構造の情報を読み出し、2403から2411の処理を順に繰り返すことで、表示用の文書の出力処理を行なう。まず、表示対象となる構

造情報を読み出す(2403)。ここで表示対象の構造であるか否かの判定は、図7を用いて前述した方法を用いて実現する。表示対象の構造情報であれば、まず、構造の開始タグを出力する(2404)。さらに本構造に下位構造が存在するなら(2405)、下位構造に対して、表示処理(2403から2411の処理)を行なう(2406)。下位構造がなくなれば、構造の終わりを示すタグを出力する処理(2411)に移る。2405のステップで下位構造が存在しないと判定された場合は、文字列の構造であるため、本構造の内容に対して、ヒット範囲情報と比較し(2408)、ヒット範囲が含まれる構造であれば、ハイライト処理を行なう(2409)。ハイライト処理については、図15を用いて前述した方法を用いる。ハイライト範囲が含まれない文字列であれば、内容をそのままテキストとして出力する(2410)。出力内容が文字列の場合は、2404、2411のステップで、開始タグ、終了タグは出力しない。上記の処理で構造ごとのハイライト表示を実現する。さらに処理すべき構造があれば、2403からの処理を繰り返す(2412)。

【0049】図25は、部分構造表示用のDTDの作成内容である。部分構造の出力により、元のDTD(2501)で必ず出現しなければならないと定義されている構造が出力されない場合がある。さらに上位構造が必ずしも出力されるとは限らない。このため、部分構造表示用のDTDは、上位構造の開始タグ、終了タグの出現を必須としない。さらに構造そのものについても、必ずしも出現しなくて良いとするように変更する必要がある。作成された部分構造表示用のDTDは2502に示したようになる。このDTDを用いて作成したSGML文書は、2503に示したようになる。この例では、タイトルだけを抽出している。

【0050】図26は、部分構造表示用のDTD作成手順を示したフローチャートである。まず、登録用のDTDを取得する(2601)。次にDTD中のELEMENT項目を取り出す(2602)。内容モデルにCDATA、RCDATA、#PCDATAが含まれる場合は、ハイライト情報を付加する(2603-2606)。ハイライト情報の付加は、図11の1103から1106ステップの処理と同じである。次に内容モデル中の出現指示子(*、+、?、なし)をチェックし、"+"ならば(2607)、"*"に変更し(2608)、出現指示子がなければ(2609)、"??"を付加する(2610)。全てのELEMENT宣言に対する処理が終わると(2611)、ハイライト用の構造のELEMENT宣言を追加し(2612)、さらに、下位構造が存在する構造のタグが出現することが必須(-)であれば、不要(o)に変更する。

【0051】次に、本発明を用いた実施例4について、図面を用いて説明する。図27は、本実施例のシステム構成図である。WWW(World Wide Web)検索シス

10

20

30

40

50

テム(2700)は、ネットワーク(2702)を使用してクライアント(2701)と接続されている。クライアント(2701)は、PC、WSなどであり、クライアント(2701)上で動作するWebブラウザ(2703)上の、検索ターム設定画面上で検索タームを入力する。WWW検索システム(2700)では、この検索タームを用いて検索を行い、その検索結果をWebブラウザ(2703)に出力する。WWW検索システム(2700)は、クライアント(2701)からの検索タームを受け取るHTTPサーバ(2704)と、検索処理およびハイライト用タグを挿入するデータ制御部(2705)と、ハイライトタグの位置情報などを格納しておくメモリ(2706)から成り立ち、検索対象となるHTML文書を格納しておく磁気ディスク装置(2707)が接続されている。データ制御部(2705)では、HTTPサーバ(2704)で受け取った検索タームを磁気ディスク(2707)中に存在するHTML文書に対して検索処理を行い、検索タームにヒットしたHTML文書の検索ヒット位置にハイライトタグを挿入する。メモリ(2706)は、各文書ごとの検索ヒット数を格納するハイライト数格納領域(2708)と、検索結果位置情報を格納するハイライト位置情報格納領域(2709)と、挿入するハイライト用タグの内容を格納しておくハイライト用タグ文字格納領域(2710)と、ハイライト用タグを挿入したHTML文書を格納するHTML文書一時格納領域(2711)と、クライアント(2701)で入力した検索タームをWWW検索システム(2700)のHTTPサーバ(2704)で取得し、一時的に格納する検索ターム格納領域(2712)からなる。WWW検索システム(2700)によってハイラ

【0052】次に、データ制御部(2705)の処理内容について、図28を用いて説明する。ここでは、クライアント(2701)で設定した検索タームを取得し、検索処理を行い、検索ヒット位置を検出しハイライト位置情報(2709)を作成し、検索条件にヒットしたHTML文書の検索タームにヒットしたHTML文書の検索ヒット位置にハイライト用のタグを埋め込み、クライアント(2701)のWebブラウザ(2703)に表示する。

ステップ2800: クライアント(2701)で設定した検索タームを、WWW検索システム(2700)では、HTTPサーバ(2704)を用いて取得する。取得した検索タームは、メモリ(2706)の検索ターム格納領域(2712)に格納される。

ステップ2801: ステップ2800で検索ターム格納領域(2712)に格納した検索タームを用いて、磁気

ディスク装置(2707)に格納されているHTML文書に対する全文検索を行う。検索ヒットした場合は、HTML文書中の検索ヒット位置や検索ヒット数などを取得し、その情報をハイライト位置情報格納領域(2709)、ハイライト数格納領域(2708)に格納する。この処理については、図29を用いて詳しく説明する。

ステップ2802: ステップ2801において、作成されたハイライト位置情報格納領域(2709)に格納されている情報を基に、ハイライトタグ文字格納領域(2710)に格納されているハイライト用タグをHTML文書の検索ヒットした位置に挿入し、HTML文書一時格納領域(2711)に格納する。詳細は、図33を用いて説明する。

ステップ2803: ステップ2802により作成されたHTML文書一時格納領域(2711)に格納されたハイライト用HTML文書を、HTTPサーバ(2704)を用いてクライアント(2701)のWebブラウザ(2703)に表示する。ステップ2800からステップ2803の処理を繰り返すことにより、クライアント(2701)で入力された検索条件を用いて、磁気ディスク(2707)に格納されているHTML文書を検索し、検索条件にヒットした文書に対して、複数箇所の検索ヒット位置のハイライト表示を可能とする。

【0053】次に、図29を用いて、図28のステップ2801のハイライト位置情報の作成処理について説明する。

ステップ2900: 磁気ディスク(2707)に格納されているHTML文書を読み出す。図34のHTML文書(3400)は、読み出したHTML文書の例である。このHTML文書をWebブラウザで表示すると、3401に示すような画面が表示される。

ステップ2901: ハイライト位置情報を格納する領域であるハイライト位置情報格納領域(2709)を α 件数分確保する。 α は、任意の正の整数である。またハイライト数を格納するハイライト数格納領域(2708)を確保する。なお、ハイライト位置情報格納領域(2709)と、ハイライト数格納領域(2708)のデータ形式は、図30および図31に示す。ハイライト位置情報格納領域(2709)は、図30に示すように、HTML文書番号(3000)、先頭からのハイライト位置番号(3001)、ハイライトバイト数(3002)、ハイライト挿入タグ番号(3003)から構成される。HTML文書番号(3000)は、ステップ2900で読み出したHTML文書の番号である。HTML文書を格納した際に付けられる通し番号などを格納する。先頭からのハイライト位置番号(3001)は、ステップ2900で読み出したHTML文書にステップ2800で取得した検索タームにヒットした場合、HTML文書中の検索ヒット位置を文書先頭からバイト数で格納する。ハイライトバイト数(3002)は、ハイライトする長

さをバイト 数で格納する。つまり、検索タームの文字列長を格納する。ハイライト 挿入タグ番号(3 0 0 3) は、複数の検索タームでハイライト 表示する場合、検索タームごとにハイライト 用タグを区別して表示することが可能である。ここに格納されている情報を基にして、ハイライト 用タグを区別する。つまり、ここには、ハイライト 表示に利用するタグの種類を判別するデータを格納する。

【 0 0 5 4 】ステップ2 9 0 2 : ハイライト 位置情報格納領域(2 7 0 9) に格納したカウントを示す *i_cnt* を 0 に初期設定する。

ステップ2 9 0 3 : ステップ2 8 0 0 で読み出した検索タームとステップ2 9 0 0 で読み出したHTML 文書が一致するか否かをチェックをする。検索ヒット 箇所が存在する場合は、ステップ2 9 0 4 に進む。また、存在しない場合は、ステップ2 9 0 8 に進む。

ステップ2 9 0 4 : ステップ2 9 0 1 または2 9 0 5 で確保したハイライト 位置情報格納領域(2 7 0 9) がハイライト 格納数を示す *i_cnt* より大きいかなをチェックする。データを格納する領域がまだ存在する場合、ステップ2 9 0 6 に進む。また、格納する領域が存在しない場合、ステップ2 9 0 5 に進む。

ステップ2 9 0 5 : ハイライト 位置情報格納領域(2 7 0 9) を一定値拡大して再度確保し直し、ステップ2 9 0 6 に進む。

【 0 0 5 5 】ステップ2 9 0 6 : ステップ2 9 0 1 または2 9 0 5 で確保したハイライト 位置情報格納領域(2 7 0 9) の *i_cnt* 番目の位置に、HTML 文書番号(3 0 0 0)、HTML 文書の先頭からの位置(3 0 0 1)、ハイライト 文字数(3 0 0 2)、ハイライトタグ挿入番号(3 0 0 3) を格納する。 *i_cnt* は0 に初期化されているので、 *i_cnt* が0 の場合、0 番目にデータを格納する。1 つのHTML 文書中に複数のハイライト 情報を格納する場合は、 *i_cnt* が更新されるので、 *i_cnt* が示す位置に格納する。ステップ2 9 0 0 で読み出したHTML 文書(3 4 0 0) をHTML 文書番号「 0 0 1 」とする。さらに、ステップ2 8 0 0 で抽出した検索タームを「 特集」とする。このHTML 文書(3 4 0 0) で、検索ターム「 特集」を検索すると、HTML 文書(3 4 0 0) の先頭から1 2 2 バイト 目(3 4 0 3) に「 特集」の文字を見つけることができる。この場合、HTML 文書番号(3 0 0 0) にはHTML 文書番号である「 0 0 1 」(3 4 0 4) を格納し、HTML 文書の先頭からの位置(3 0 0 1) には「 1 2 2 」(3 4 0 5) を格納し、ハイライト 文字数(3 0 0 2) には「 特集」のバイト 数「 4 」(3 4 0 6) を格納する。最後に、ハイライトタグ挿入番号(3 0 0 3) には、検索結果を強調するためのタグを示す番号を格納する。ここでは、「 1 」(3 4 0 7) を格納する。

【 0 0 5 6 】ここで、ハイライト 挿入タグ番号と実際に

格納するハイライトタグを対応する構成を図3 2 に示す。図3 2 の(1) では、ハイライトタグ文字格納領域(2 7 1 0) に格納されているハイライト 挿入タグ用の構造体3 2 0 0 を示す。ハイライト 挿入タグ用の構造体(3 2 0 0) は、通し 番号を格納するタグ番号1 (3 2 0 2) と、ハイライト 開始タグ名を格納する開始タグ1 (3 2 0 3)、ハイライト 終了タグ名を格納する終了タグ1 (3 2 0 4) と、タグの個数を格納するハイライトタグ数(3 2 0 1) から成り 立つ。ハイライトタグ数に格納した数分のタグ番号、開始タグ、終了タグが存在する。

【 0 0 5 7 】ハイライトタグ文字格納領域の使用例を(2) に説明する。ここでは、3 種類のハイライトタグを格納する場合を示す。よって、ハイライトタグ数を格納する箇所には、「 3 」(3 2 0 5) を格納する。タグ番号「 0 」(3 2 0 6) の箇所の開始タグには、赤色を示すタグ「 」(3 2 0 7) を、終了タグには「 」(3 2 0 8) を格納する。同様に、タグ番号「 1 」(3 2 0 9) には、点滅を示すタグ「 <BLINK > 」を、タグ番号「 2 」(3 2 1 0) には、文字を大きく表示する「 <H1 > 」を格納する。ハイライトタグ文字格納領域(2 7 1 0) は、ハイライト 位置情報格納領域(2 7 0 9) の作成前に作成する。また、このハイライトタグ文字格納領域(2 7 1 0) は、ユーザインターフェースを使用して、作成することも可能である。複数のハイライト 用タグを用意することで、異表記や同義語の検索処理を行った場合において、異表記で検索された文字にはタグ番号「 1 」、同義語で検索された文字にはタグ番号「 2 」のように、検索条件ごとに異なるハイライト 表示が可能となる。ハイライト 用タグに「 <BLINK > 」を使用する場合は、ハイライト 位置情報格納領域(3 4 0 2) のハイライトタグ挿入番号(3 4 0 7) に「 1 」を格納する。

【 0 0 5 8 】ステップ2 9 0 7 : ステップ2 9 0 6 において、ハイライト 位置情報格納領域(2 7 0 9) にデータを格納したので、 *i_cnt* を1 を加え、ステップ2 9 0 3 に戻る。

ステップ2 9 0 8 : ステップ2 9 0 0 で取得したHTML 文書中のハイライト 数を取得し、ハイライト 数格納領域(2 7 0 8) に格納する。ハイライト 数格納領域(2 7 0 8) の構造体の内容は図3 1 を用いて説明する。図3 1 は、ハイライト 数格納領域(2 7 0 8) の構造体の内容である。3 1 0 0 は、ステップ2 9 0 0 で読み出したHTML 文書の文書番号である。また、3 1 0 1 は、取得したハイライト 数を格納しておく 箇所である。ここでは、文書番号「 0 0 1 」を文書番号3 1 0 0 に格納し、 *i_cnt* をハイライト 数格納領域(3 1 0 1) に格納し、処理を終了する。

【 0 0 5 9 】次に、図3 3 を用いて、ハイライト 用タグ

付のHTML文書作成処理について説明する。

ステップ3300: ステップ2900で読み出したHTML文書中に、ハイライトタグを挿入する必要があるかをチェックする。ハイライト位置情報格納領域(2709)に格納したHTML文書番号(3000)が存在する場合は、ステップ3301に進む。存在しない場合は、ステップ3309ですべてのテキストを出力し、処理を終了する。

ステップ3301: 処理カウントを示すi_cntを0に初期化する。

ステップ3302: ハイライトタグを挿入したHTML文書を格納するHTML文書一時格納領域(2711)を確保する。HTML文書一時格納領域(2711)は、HTML原文書のバイト数は、ハイライト用開始タグと終了タグのバイト数の合計値にハイライト挿入数を乗じたバイト数の領域を確保する。ハイライトの開始タグと終了タグは、ハイライト位置情報格納領域(2709)のハイライト挿入タグ番号(3003)より、ハイライト用タグのタグ文字列長を計算する。ハイライト数は、ステップ2908でハイライト数格納領域(2708)に格納したハイライト数(3101)を取得する。

ステップ3303: ハイライト数(3101)がi_cntより小さいか否かをチェックする。小さい場合は、未処理のハイライト箇所が存在するので、ステップ3304に進む。それ以外は、処理すべき未処理のハイライト箇所を全て終了したので、ステップ3309に進む。

ステップ3304: ハイライト位置までのHTML文書をステップ3302で確保したHTML文書一時格納領域(2711)に格納する。
【0060】ステップ3305: ハイライト開始タグをHTML文書一時格納領域(2711)に格納する。ハイライト開始タグは、ハイライト挿入タグ番号(3003)から抽出した番号より得られるハイライトタグ文字格納領域(2710)に格納されているタグ名である。図34(3)の場合、ハイライト挿入タグ番号(3003)には「1」が格納されている。図32(2)に示したハイライトタグ文字格納領域(2710)のタグ番号「1」(3209)に格納されている「<BLINK>」をHTML文書一時格納領域(2711)に格納する。

ステップ3306: 検索タームをHTML文書一時格納領域(2711)に格納する。図34の場合、「特集」をHTML文書一時格納領域(2711)に格納する。ステップ3307: ハイライト終了タグをHTML文書一時格納領域(2711)に格納する。ハイライト終了タグは、ステップ3305で処理したハイライト開始タグ同様、ハイライト挿入タグ番号(3003)にて格納された番号から得られるハイライトタグ文字格納領域(2710)に格納されているタグ名を格納する。図34(3)の場合、「1」が格納されている。したがっ

て、図32(2)のタグ番号「1」に格納されている「</BLINK>」をHTML文書一時格納領域(2711)に格納する。

ステップ3308: ステップ3305からステップ3307において、データをHTML文書一時格納領域(2711)に格納した後、i_cntに1を加え、ステップ3303に戻る。

ステップ3309: ハイライト挿入位置からHTML文書最後までテキストをHTML文書一時格納領域(2711)に格納し、ハイライトタグ付きHTML文書の作成処理を終了する。

【0061】以上の処理を用いることで、クライアント(2701)設定した検索タームから、HTML文書を検索し、検索タームと一致する文書に対して、ハイライト数を格納するハイライト数格納領域(2708)、ハイライト位置を格納するハイライト位置情報格納領域(2709)の内容を作成することが可能である。上記の処理結果の例を図35に示す。3500は、ハイライト用タグを挿入したHTML文書である。検索ヒットした「特集」の前後(3501, 3502)にハイライト用タグが挿入されている。このHTML文書を画面に表示すると3503のようになり、検索ヒットした「特集」(3504)が点滅表示される。以上で、本発明の第1実施例として、クライアント(2701)が挿入した検索タームをHTML文書(2707)の中から検索し、ヒット位置にハイライト用タグを挿入した、ハイライト用タグ付きHTML文書を作成する方法を説明した。

【0062】次に、本発明における実施例5について、図36から図42を用いて説明する。図36は、検索タームがHTML文書のタグで分断されている場合や、検索タームがタグ内に存在する場合のハイライト表示方法におけるシステム構成図である。図27と同様に、クライアント(2701)のWebブラウザ(2703)上で検索タームを設定する。

【0063】WWW検索システム(2700)は、検索タームを取得するHTTPサーバ(2704)、検索処理を行うデータの制御(2705)、領域を確保するメモリ(2706)から成り立つ。メモリ(2706)は、図27の説明で述べた以外に、レイアウト表示などに使用されるタグで、読み飛ばすタグ名を格納した読み飛ばしタグ名格納領域(3600)と、クライアント(2701)が入力した検索タームとHTML文書(2707)が一致した開始位置を一時的に格納しておく開始位置格納領域(3601)と、検索タームとヒットした位置がHTMLタグの開始文字「<」と終了文字「>」の間に存在した場合、目印となるマークを格納しておく再度記述マーク格納領域(3602)と、HTMLタグの開始タグと終了タグの間に検索タームがヒットした場合、検索ヒットした箇所の前後にハイライト用タ

グを入れることができないHTMLタグを記述しておく、ハイライトタグ挿入不可能タグ名格納領域(3603)からなる。検索タームがHTML文書のタグをまたがっている場合や、検索タームがタグ内に存在する場合の検索タームの取得、ハイライト位置情報の作成、ハイライト用タグ挿入方法は、図28で示した処理手順で行う。また、各々の処理内容については、図37から図42を用いて説明する。

【0064】ステップ2800で取得した検索タームを用いて、ステップ2801の処理では、検索処理およびハイライト位置情報作成処理を行う。処理内容は、図37のフローチャートに示す。

ステップ3700：処理対象となるHTML文書を磁気ディスク(2707)から読み出す。

ステップ3701：ハイライト位置情報を格納するハイライト位置情報格納領域(2709)とハイライト数格納領域(2708)をメモリ(2706)に確保する。

ステップ3702：検索ヒット位置の前後に挿入するハイライトタグを読み出す。図32(2)の使用例に具体例を示したようにハイライトタグ文字格納領域(2710)からハイライト用タグを読み出す。この場合ハイライト挿入タグ番号の個数は、「3」(3205)から「3つ」とわかる。1番目の「0」(3206)には、「」(3207)と「」(3208)格納されている。そこで、ハイライト挿入タグ番号0番目の開始タグは「」、終了タグは「」となる。同様に、ハイライト挿入タグ番号1番目の開始タグは「<BLINK>」、終了タグは「</BLINK>」となり、ハイライト挿入タグ番号2番目の開始タグは「<H1>」、終了タグは「</H1>」となる。

ステップ3703：HTML文書の処理済み文字数のカウントを示す i_cnt と、ハイライト数を格納する領域の内容を0に初期設定する。

【0065】ステップ3704：検索タームとHTML文書の文字列が一致するか否かをチェックする。チェック方法として、HTML文書の i_cnt バイト目から、検索タームの先頭文字と一致する文字を検索する。ステップ3703において、初期設定が0に設定されているため、最初は、HTML文書の0バイト目から一致する文字を検索する。一致した場合は、ステップ3705に進む。不一致の場合は、処理を終了する。また、ここでは、検索タームを抽出する方法として、指定したタグを飛ばして検索する方法を用いる。具体的には、読み飛ばしタグ名格納領域(3600)に格納してあるタグ名をHTML文書中に存在した場合は、そのタグを読み飛ばし、検索処理を行う。読み飛ばしタグ名格納領域(3600)に「IMG」を格納しておき、図34のHTML文書(3400)を検索した場合、HTML文書(34

00)中の先頭からデータを走査し、「IMG」(3413)が抽出された時点で、タグ内の文字を飛ばす。つまり、タグの終了文字「>」(3414)までを飛ばす。この読み飛ばしタグ名格納領域(3600)は、検索処理前に設定しておくことにより、読み飛ばし処理が可能となる。

【0066】ステップ3705：ステップ3704でHTML文書の先頭から検索タームの先頭文字と一致した文字までの文字数を開始位置格納領域(3601)に一時的に確保する。

ステップ3706：検索タームの文字列とHTML文書に書かれている文字が一致するか否かをチェックし、一致した場合、一致箇所がHTMLタグ内に存在するかあるいはHTMLタグ外に存在するか否かをチェックする。さらに、検索ヒットした文字列の最後の文字の位置を、HTML文書の先頭からの文字数で確保する。詳細は、図38を用いて説明する。

ステップ3707：ステップ3706の結果、検索ヒットしたか否かをチェックする。HTML文書中に検索タームが存在した場合は、ステップ3708に進む。検索タームが存在しない場合、ステップ3712に進む。

ステップ3708：ステップ3701で確保したハイライト数格納領域(3708)とハイライト格納数を比較して、確保した領域が格納したハイライト数より多ければ、ステップ3709に進む。少なければ、ステップ3710に進む。

ステップ3709：ハイライト位置情報格納領域(2709)にデータを格納する領域が足りないため、再度領域設定し直し、ステップ3710に進む。

【0067】ステップ3710：ハイライトする文字数とハイライトの位置の情報を、ハイライト位置情報格納領域(3600)に格納する。具体的には、図30で説明したハイライト位置情報格納領域(3600)のHTML文書番号(3000)には、ステップ3700で読み出したHTML文書の文書番号を格納し、先頭からのハイライト位置情報(3001)には、ステップ3705で取得した開始位置を格納する。また、ハイライトのバイト数(3002)には、検索タームの文字列長を格納し、ハイライト挿入タグ番号(3003)には、ステップ3702で読み出したタグの番号を格納する。ハイライト挿入タグ番号(3003)は、デフォルトとして、「0」を設定する。

ステップ3711：検索タームにヒットする文字列が複数存在する場合、検索ヒットした位置の次文字から再度検索タームとHTML文書中の一致する箇所をチェックする処理を行う。そこで、ステップ3706で確保した検索ヒットの最後の文字が記述されている位置の、HTML文書の先頭からの文字数に1を加えた値を i_cnt に代入する。処理位置を更新したら、ステップ3704に戻る。

ステップ3712: ステップ3705で取得した開始位置格納領域(3600)に格納してある開始位置からの文字列と、検索タームが一致していない場合、開始位置の次文字から再度検索タームとHTML文書中の一致する箇所をチェックする処理を行う。そこで開始位置格納領域(3600)に格納してある開始位置に1を加えた値を*i_cnt*に代入する。処理位置を更新したら、ステップ3704に戻る。以上で、タグ内およびタグ外のチェックを含む検索処理およびハイライト位置情報作成処理について述べた。

【0068】次に、図38を用いて、ステップ3706のタグ内の検索およびタグ外の検索処理について説明する。ここでは、ステップ3705で取得した検索ヒットの開始位置が、文書構造を示すタグの属性中に存在するかあるいはタグの外に存在するかをチェックし、検索ヒットの開始位置からの文字列が検索タームと一致するか否かのチェックを行う。

ステップ3800: ステップ3705で開始位置格納領域(3600)に格納した検索ヒットの開始位置において、HTMLタグ内かあるいはタグ外かをチェックする。ステップ3706時点におけるHTML文書の*i_cnt*バイト目から、検索ヒットの開始位置までのデータをチェックする。タグの開始文字「<」と対応するタグの終了文字「>」をチェックし、タグ内に検索ヒットの開始位置が存在するか否かをチェックする。タグの開始文字「<」があり、タグの終了文字「>」の前に検索ヒットの開始位置が存在する場合、開始位置はタグ内に存在するとして、ステップ3801に進む。タグの開始文字「<」とタグの終了文字「>」に囲まれない範囲に、検索ヒットの開始位置が存在する場合、検索ヒットの開始位置は、タグ外に存在するとして、ステップ3804に進む。

【0069】ステップ3801: 検索タームと、検索ヒットの開始位置からの文字列が一致するか否かをチェックする。検索タームの文字列が複数バイトから成り立つ場合、文字列を1バイト毎にチェックする。検索タームの文字列と検索ヒットした位置からの文字列が一致する場合、ステップ3802に進む。不一致の場合、ステップ3803に進む。

ステップ3802: ステップ3801において、検索タームと一致した場合、「検索ヒット」として、処理を終了する。また、検索ヒットした文字列の終端位置を求める。終端位置は、検索ヒットした開始文字位置に検索タームの文字列長を加えたバイト数とする。ここで求めた終端位置は、ステップ3711にて使用される。

ステップ3803: ステップ3801において、検索タームが不一致の場合、「検索ヒットしない」として、処理を終了する。

ステップ3804: ステップ3800において、検索ヒットの開始位置がタグの外に存在した場合、タグ外用の

検索処理を行う。タグ外用の検索処理は、図39を用いて説明する。

ステップ3805: ステップ3804で検索タームがヒットする箇所がHTML文書中に存在するか否かをチェックする。存在する場合は、ステップ3807に進む。存在しない場合は、ステップ3806に進む。

ステップ3806: ステップ3805において、検索タームがヒットしない場合、処理を終了する。

ステップ3807: ステップ3805において、検索タームと一致した場合、「検索ヒット」として、処理を終了する。また、検索ヒットした文字列の終端位置を求める。終端位置は、検索ヒットの開始文字位置に、ステップ3804で検出した検索ヒットの最後の文字が記述されている位置を加えた値とする。ここで求めた終端位置は、ステップ3711にて使用する。以上で、タグ内検索およびタグ外検索処理について説明した。

【0070】次に、ステップ3804のタグ外用の検索処理について、図39を用いて説明する。

ステップ3900: HTML文書中に検索タームが存在するか否かをチェックする。検索タームの文字列がHTML文書中に存在する文字列と一致するか否かをチェックするが、検索ヒットした開始位置から、途中に存在するタグを飛ばすことにより一致する場合があるので、ここでは、検索ヒットした開始位置から1文字ごとに検索タームと合致しているか否かをチェックする。具体的には、図34を用いて説明する。検索タームを「特集記事」とした場合、(2)の表示画面では、3408に「特集記事」が表示されている。しかし、HTML文書(3400)では、「特集」(3403)と「記事」(3416)の間に「</H1>」(3417)のタグがある。このように検索タームの途中にHTMLタグが存在する場合、このHTMLタグを読み飛ばして、検索タームと一致する文字列を抽出する。ここでは、検索タームを1文字ごととHTML文書の文字と照合し、チェックを行う。検索タームの1文字目とHTML文書中の文字が一致した場合は、検索タームの次文字とHTML文書の次文字について、同処理を繰り返す。検索タームの文字列のすべての文字が一致した場合、具体的には、「特」(3403)、「集」と文字の比較を行い、「</H1>」(3417)を読み飛ばし、さらに、「記」(3416)、「事」と文字比較を行う。すべての検索タームを抽出し終わった場合、ステップ3901に進む。検索タームがHTML文書中の文字列と完全に一致しなかった場合、ステップ3902に進む。

【0071】ステップ3901: HTML文書中に検索タームが存在するため、「検索ヒット」として、処理を終了する。また、検索ヒットの終端位置を求める。終端位置は、ステップ3900において、最後に抽出した文字の位置である。

ステップ3902: ステップ3900で、検索タームの

文字とHTML文書の文字が一致なかった場合、HTML文書の文字が、タグの開始文字「<」か否かをチェックする。タグの開始文字「<」の場合は、ステップ3903に進む。それ以外の文字の場合は、ステップ3904に進む。

ステップ3903: ステップ3902において、HTML文書中の文字がタグの開始文字「<」の場合、タグの内容を飛ばして、ステップ3900に戻る。具体的には、タグの終了文字「>」を抽出し、抽出した文字までを読み飛ばす。図34のHTML文書(3400)で、検索タームを「特集記事」とした場合、「特集」(3403)の次文字にある「<」(3417)からタグの終了文字「>」(3418)までを読み飛ばす。つまり、「</H1>」を読み飛ばす。

ステップ3904: ステップ3902において、検索タームが不一致の場合、「検索ヒットしない」として、処理を終了する。以上で、HTML文書中に検索タームの文字列が存在した場合、検索ヒット位置を抽出し、ハイライト位置情報格納領域にハイライト位置情報を格納する処理について説明した。

【0072】次に、図40を用いて、ハイライト位置情報格納領域に格納したハイライト位置情報を基にして、HTML文書の検索ヒットした文字列を強調表示するため、ハイライト用タグを挿入する方法について説明する。

ステップ4000: ステップ3710においてハイライト情報格納領域(2709)に格納した、ハイライト位置情報を読み出す。

ステップ4001: ハイライトタグを挿入したHTML文書を格納するためのHTML文書一時格納領域(2711)を確保する。確保する領域の大きさは、HTML文書の原文書のデータ、ハイライトタグ数分のハイライト開始タグと終了タグの長さの和を乗じた値のバイト数である。ハイライトタグ数は、ハイライトタグ数格納領域(2708)から読み出す。また、ハイライトの開始タグと終了タグは、ハイライト位置情報格納領域(2709)のハイライト挿入タグ番号(3003)とハイライトタグ文字格納領域(2710)からタグを検出し、検出したタグの文字列長を求める。

ステップ4002: HTML文書中の処理済み位置を示すi_cntと、ハイライト処理数を0に初期化する。

【0073】ステップ4003: 処理済みのハイライト箇所の数であるハイライト処理数と、処理すべきハイライト数を比較する。ハイライト処理数が少ない場合は、ハイライト用タグを挿入する処理を行うため、ステップ4004に進む。それ以外の場合は、ステップ4007に進む。

ステップ4004: 処理済みの位置を示すi_cntから検索ヒットの開始位置までのデータを、HTML文書一時格納領域(2711)に格納する。具体的には、図34

のHTML文書(3400)で、検索タームを「特集記事」とした場合、HTML文書先頭から「特集記事」(3403)前の文字「今月の」までのデータをHTML文書一時格納領域(2711)に格納する。

ステップ4005: ハイライト用タグを検索ヒット位置に格納する。ハイライト用タグの挿入処理については、図41で説明する。

ステップ4006: HTML文書の処理済み位置を示すi_cntにハイライト終了タグを挿入した位置の先頭からのバイト数を代入し、ステップ4003に戻る。

ステップ4007: HTML文書の処理済み位置を示すi_cntから、HTML文書の最後までまでのデータをHTML文書一時格納領域(2711)に格納し、処理を終了する。

【0074】次にステップ4005で処理するハイライトタグの挿入処理について、図41を用いて説明する。ここでは、検索ヒットした位置が、タグの内あるいはタグ外かをチェックし、検索ヒット位置の前後にハイライト用タグを挿入する処理を行う。

ステップ4100: HTML文書でヒットした位置がHTMLタグ内か、タグ外かをチェックする。チェック方法は、検索ヒットの開始位置までのHTML文書において、HTMLタグの開始文字「<」とタグの終了文字「>」の対応をとり、タグ内かタグ外かを判断する。検索ヒットの開始位置がタグの開始文字「<」からタグの終了文字「>」の間にある場合は、タグ内に存在するとし、ステップ4101に進む。それ以外の場合は、タグ外に検索ヒット位置の開始位置が存在するとし、ステップ4110に進む。

ステップ4101: タグの開始文字「<」の次文字から文字を抽出し、タグの種類を取得する。例えば、図34のHTML文書(3400)の場合、検索タームを「hitachi」とした場合、HTML文書(3400)中の3409に「hitachi」を取得することができる。このHTMLタグの種類を取得すると、タグの開始文字「<」の次に書かれている「A」(3410)とわかる。

ステップ4102: ステップ4101で取得したタグが、開始用のタグが終了用のタグかをチェックする。終了用のタグの場合、タグの開始文字「<」の次文字が「/(スラッシュ)」である。そこで、タグの開始文字「<」の次文字をチェックし、判別する。このタグ開始文字「<」の次文字が「/」の場合、終了用のタグと判定して、ステップ4105に進む。それ以外の場合は、開始タグと判定し、ステップ4103に進む。

【0075】ステップ4103: 開始用タグと終了用タグの間にハイライト用のタグを挿入することが可能かをチェックする。挿入することが可能な場合は、ステップ4105に進む。また、不可能な場合は、ステップ4106に進む。具体的には、図34のHTML文書(3400)で、検索タームが「hitachi」の場合、H

20

30

40

50

ステップ4201: 検索ヒットした文字列の前後にハイライト用タグを挿入することが可能か否かをチェックする。チェック方法は、検索ヒットした検索文字列に囲まれているHTMLタグを抽出する。抽出したタグの種類とハイライトタグ挿入不可能タグ名格納領域(3603)に格納されているタグと比較する。一致すれば、ステップ4202に進み、不一致の場合、ステップ420

9に進む。ハイライトタグ挿入不可能タグ名格納領域(3603)中に記述されたタグ名は、開始用タグと終了用タグの間に、ハイライト用タグを挿入することができない。よって、ハイライトタグ挿入不可能タグ名格納領域(3603)に格納されているHTMLタグと比較し、一致すれば、ステップ4202に進む、不一致の場合は、ステップ4209に進む。このハイライトタグ挿入不可能タグ名格納領域(3603)は、ユーザインタフェースを使用して、データ制御(2705)の前に作成しておく。

ステップ4202: ハイライト用タグを挿入することが出来ない場合、終了用タグの終わりの文字「>」までのHTML文書を読み飛ばす。図34において、検索ターム「日立」とした場合、「」(3412)までのHTML文書を読み飛ばす。

【0078】ステップ4203: ステップ4202で飛ばしたHTML文書をHTML文書一時格納領域(2711)に格納する。図34において、検索ターム「日立」とした場合、ステップ4002で設定したi_cnt番目あるいはステップ4006で更新したi_cnt番目にあるHTML文書の文字から「」(3412)のデータをHTML文書一時格納領域(2711)に格納する。

ステップ4204: ハイライト用タグの開始タグをHTML文書一時格納領域(2711)に挿入する。ハイライト位置情報格納領域を3402として、ハイライトタグ文字格納領域を図32の(2)とした場合、「<BLINK>」が抽出される。よって、ここでは、「<BLINK>」を挿入する。

ステップ4205: 再表示用マークを格納する。ステップ4108同様に、再度記述マークの格納領域(3602)に格納されているHTML文書を読み出し、HTML文書一時格納領域(2711)に格納する。

ステップ4206: 検索ヒットした文字列をもう一度HTML文書一時格納領域(2711)に挿入する。但し、検索ヒットした文字列の途中にタグが存在する場合は、タグが存在する箇所までの文字列を挿入する。

ステップ4207: ハイライト用の終了タグをHTML文書一時格納領域(2711)に挿入する。ここでは、「</BLINK>」を挿入する。

ステップ4208: ステップ4206において、検索タームの文字列をすべて格納したか否かをチェックする。検索ヒットした文字列中にタグが存在し、検索ヒットした文字をすべて格納していない場合、ステップ4200に戻る。また、すべての文字を格納した場合は、処理を終了する。

【0079】ステップ4209: 検索ヒットした開始位置までデータを飛ばし、飛ばしたHTML文書をHTML文書一時格納領域(2711)に格納する。具体的には、図34のHTML文書(3400)で、検索ターム

が「特集記事」とした場合、検索ヒットした「特集」(3403)の前に存在する「今月の」までのHTML文書をHTML文書一時格納領域(2711)に挿入する。

ステップ4210: ステップ4205同様に、ハイライト用タグの開始タグをHTML文書一時格納領域(2711)に格納する。ここでは、「<BLINK>」を挿入する。

ステップ4211: 検索ヒットした文字列を、HTML文書一時格納領域(2711)に挿入する。但し、検索ヒットした文字列の途中にタグが存在する場合は、タグが存在する箇所までの文字列を挿入する。例えば、HTML文書(3400)で、検索タームが「特集記事」とした場合、「特集」(3403)と「記事」(3416)の間に「</H1>」(3417)が存在する。よって、ここでは、「特集」を格納する。

ステップ4212: ハイライト用タグの終了タグをHTML文書一時格納領域(2711)に挿入する。ここでは、「</BLINK>」を挿入する。

【0080】ステップ4213: ステップ4211において、検索タームの文字列すべてをHTML文書一時格納領域(2711)に挿入したか否かをチェックする。HTMLのタグを除くことにより、検索タームとHTML文書の文字列がヒットし、検索ヒットした先頭位置から検索ターム長の文字列の間に、HTMLのタグが存在する場合、ステップ4211では、HTMLタグまでのHTML文書をHTML文書一時格納領域(2711)に挿入する。この場合、HTMLタグから残りの検索ヒットの文字を処理する必要がある。すべての検索タームをHTML文書一時格納領域に挿入した場合は、処理を終了する。また、HTMLタグから残りの検索ヒットの文字を処理する場合は、ステップ4200に戻る。図34のHTML文書(3400)で、検索タームが「特集記事」とした場合、「特集」(3403)と「記事」(3416)の間に「</H1>」(3417)が存在する。ステップ4206では、「特集」のみ挿入した状態で、「記事」を挿入していないため、ステップ4200に戻る。

【0081】このような処理を行うことで、クライアント(2701)が設定した検索タームを用いて、検索タームと合致するHTML文書にハイライトタグを挿入し、Webブラウザ(2703)にハイライトヒット箇所を表示することが可能である。ここでは、検索ターム1つに対して、HTML文書をチェックし、検索タームの文字列がHTML文書中に存在すれば、クライアント(2701)のWebブラウザに検索の結果を表示する処理方法を示したが、1つの検索タームに対して、複数のHTML文書から検索し、検索ヒットしたHTML文書数分のハイライト位置情報を格納し、ハイライト用タグを格納した複数HTML文書を一括して作成する事も可

能である。また、複数の検索タームに対して、複数のHTML文書から検索し、検索ヒットしたHTML文書数分のハイライト位置情報を格納し、ハイライト用タグを格納した複数HTML文書を一括して作成することも可能である。

【0082】次に本発明を用いた実施例6について説明する。本実施例の実施例2からの変更点は、検索条件中に検索タームなどと共に、検索条件にヒットした場合のハイライト方法を定義することができる点である。これにより、任意の検索条件に対して、検索条件毎にハイ

ライト方法を指定することができる。本実施例のシステム構成は図1と同じである。ただし、検索条件103の記載方法が異なる。本実施例における検索条件103の記述方法の例を図43を用いて説明する。

【0083】図43に本実施例における検索条件の例を示す。本図に示すように、各検索タームや構造条件などの後ろに、「{アンダーライン}」のようにハイライト方法を指定する。実施例2における検索条件は、「検索対象の構造指定: 検索条件式」であったが、「検索対象の構造条件{ハイライト方法}: ハイライト方法付き検索条件式」となる。ハイライト方法の指定は省略可能である。省略時は、実施例2で示した方法でハイライト表示を行なう。すなわち、ハイライト方法が検索条件中に記載されていない箇所については、図18に示したハイライト方法定義1801を読み出し、本定義情報に記載されているハイライト方法を用いてハイライト表示する。

【0084】図44に本実施例におけるヒット範囲情報4401の格納内容を示す。実施例2の図17に示したヒット範囲情報からの変更点は、各ヒット範囲ごとにヒット条件4402だけではなく、ハイライト方法4403を格納する点である。本情報は、図43を用いて前述した検索条件を解析し、検索条件中に記載されたハイライト方法の情報を読み出すことで取得可能である。

【0085】図45に本実施例におけるハイライト表示用DTDの生成方法を示す。本例では、検索の度に新規にハイライト方法が変更される可能性があることから、ハイライト表示の度に、必要な構造だけを追加したハイライト表示用DTDを生成することとする。この場合、DTD中に検索条件ではなく、直接ハイライト方法に関する記述を行なうことになる。本図に示すように、登録に用いた元のDTD(1901)に対して、上位のハイライト構造内には下位のハイライト構造を階層的に指定でき、さらに省略も可能なように定義を変更、追加したハイライト表示用のDTD(4501)を生成している。

【0086】DTDの作成方法は、まず図44のヒット範囲情報にハイライト方法4403が記載されていない場合に、図18に示したハイライト方法定義からヒット条件に対応するハイライト方法を取得する。まず、元のDTDの各構造に対して、下位構造に出現するハイライ

ト方法を内容モデルに持つことができるように、構造情報を変更する(4502)。さらに、ヒット範囲情報4401におけるヒット範囲の階層関係から、出現するハイライト用構造の階層関係を得る。ここで得られたハイライト表示用の階層関係を元に、各ハイライト構造の下位構造として、下位のハイライト構造および文字列を内容モデルとして持つようにする。下位のハイライト構造がなければ、内容モデルとして、文字列だけが出現するようにする(4503)。

【0087】本実施例におけるハイライト処理により、検索条件をハイライト構造とするのではなく、記載されたハイライト方法を記述したハイライト表示用構造化文書と、ハイライト表示用のDTDを生成することになる。このように、本実施例による処理により、ハイライト表示用の構造化文書が生成される。生成されたハイライト表示用の構造化文書の例を図46に示す。図46に示すハイライト表示用の構造化文書をハイライト表示すると、図47に示すようになる。

【0088】

【発明の効果】本発明により、構造化文書の検索結果として、ヒットした文書の内容を表示する際に、各構造ごとに検索タームがヒットした範囲に、ハイライト情報を付加した構造化文書を出力することが可能となる。ブラウザ依存のハイライト情報ではなく、構造化文書中にハイライト情報を埋め込むことで、どのようなブラウザにおいてもハイライト表示が可能となる。検索時の条件、または、各検索タームの重要度、出現頻度などの条件によって異なるハイライト処理が行え、重要な検索タームについては、高い重み付けであることを明示したハイライト処理を行なうことが可能となる。さらに、検索条件中にハイライト方法を記述することで、ユーザ毎に任意のハイライト表示を行なうことが可能となる。さらに、部分構造だけを抽出して、ハイライト情報を付加した構造化文書を出力することが可能になる。また、文書構造を示すHTMLタグが存在する文書から文字列を検索する場合、設定した検索タームと一致した文字列がHTMLタグ内に存在する場合や、検索タームがHTMLタグをまたがって記述されている場合でも容易に検索ヒットすることが可能となる。また、検索ヒットした文字列をハイライト表示することが可能となる。

【図面の簡単な説明】

【図1】実施例1、2の構造化文書検索表示装置の処理ブロック図である。

【図2】構造化文書検索表示処理のフローチャートを示す図である。

【図3】構造化文書登録の内容を示す図である。

【図4】構造化文書登録処理のフローチャートを示す図である。

【図5】検索用のテキストを示す図である。

【図6】更新処理のフローチャートを示す図である。

【図7】指定構造の抽出処理のフローチャートを示す図である。

【図8】構造指定の解析結果として出力される情報を示す図である。

【図9】文書表示処理のフローチャートを示す図である。

【図10】構造化文書およびハイライト処理結果の例を示す図である。

【図11】文書表示用DTD作成処理のフローチャートを示す図である。

【図12】構造化文書検索用の正規化処理の内容を示す図である。

【図13】正規化処理を行なった結果の格納内容を示す図である。

【図14】正規化処理を行なった場合のヒット範囲情報の変換処理内容を示す図である。

【図15】正規化処理を行なった場合のヒット範囲情報の変換処理のフローチャートを示す図である。

【図16】ハイライト情報を付加する処理のフローチャートを示す図である。

【図17】実施例2におけるヒット範囲情報を示す図である。

【図18】実施例2におけるヒット情報ごとのハイライト方法の定義を示す図である。

【図19】実施例2のハイライト表示用DTDへの変換内容を示す図である。

【図20】実施例2におけるハイライト処理のフローチャートを示す図である。

【図21】実施例2によりハイライト情報を付加したSGML文書の例を示す図である。

【図22】ハイライト表示の例を示す図である。

【図23】実施例3の構造化文書検索表示装置の概略処理ブロック図である。

【図24】実施例3の処理内容のフローチャートを示す図である。

【図25】部分構造表示用のDTDへの変換処理を示す図である。

【図26】部分構造表示用のDTD作成処理のフローチャートを示す図である。

【図27】実施例4におけるシステム構成図である。

【図28】データ制御部のフローチャートを示す図である。

【図29】実施例4における文字検索処理およびハイライト位置情報の作成処理のフローチャートを示す図である。

【図30】ハイライト位置情報格納領域の構成である。

【図31】ハイライト数格納領域の構成である。

【図32】ハイライトタグ文字格納領域の構成である。

【図33】実施例4におけるハイライトタグ付きHTML文書の作成処理のフローチャートを示す図である。

【図34】ハイライト挿入例である。

【図35】ハイライト挿入後の例である。

【図36】実施例5におけるシステム構成図である。

【図37】実施例5における検索処理およびハイライト位置情報作成処理2のフローチャートを示す図である。

【図38】実施例5におけるタグ内検索およびタグ外検索処理のフローチャートを示す図である。

【図39】実施例5におけるタグ外用検索処理のフローチャートを示す図である。

【図40】実施例5におけるハイライト用タグの挿入HTML文書の作成処理のフローチャートを示す図である。

【図41】実施例5におけるハイライトタグ挿入処理のフローチャートを示す図である。

【図42】実施例5におけるタグ外ハイライトタグ挿入処理のフローチャートを示す図である。

【図43】実施例6における検索条件の例である。

【図44】実施例6におけるヒット範囲情報の例である。

【図45】実施例6におけるハイライト表示用DTDへの変換処理を示す図である。

【図46】実施例6におけるハイライト表示用のSGML文書の例を示す図である。

【図47】実施例6におけるハイライト表示の例を示す図である。

【符号の説明】

101 構造化文書検索表示装置

102 登録用構造化文書

103 検索条件

104 文書登録処理モジュール

105 構造化文書DB

106 検索用情報DB

107 構造化文書読み出し処理モジュール

108 検索処理モジュール

109 ヒット文書番号情報

110 ヒット範囲情報

111 ヒット文書の文書内容

112 文書表示処理モジュール

113 表示用文書

114 登録用文書格納ファイル

115 入出力装置

2301 表示構造情報

2302 部分構造表示モジュール

2700 WWW検索システム

2701 クライアント

2703 Webブラウザ

2704 HTTPサーバ

2705 データ制御

2706 メモリ

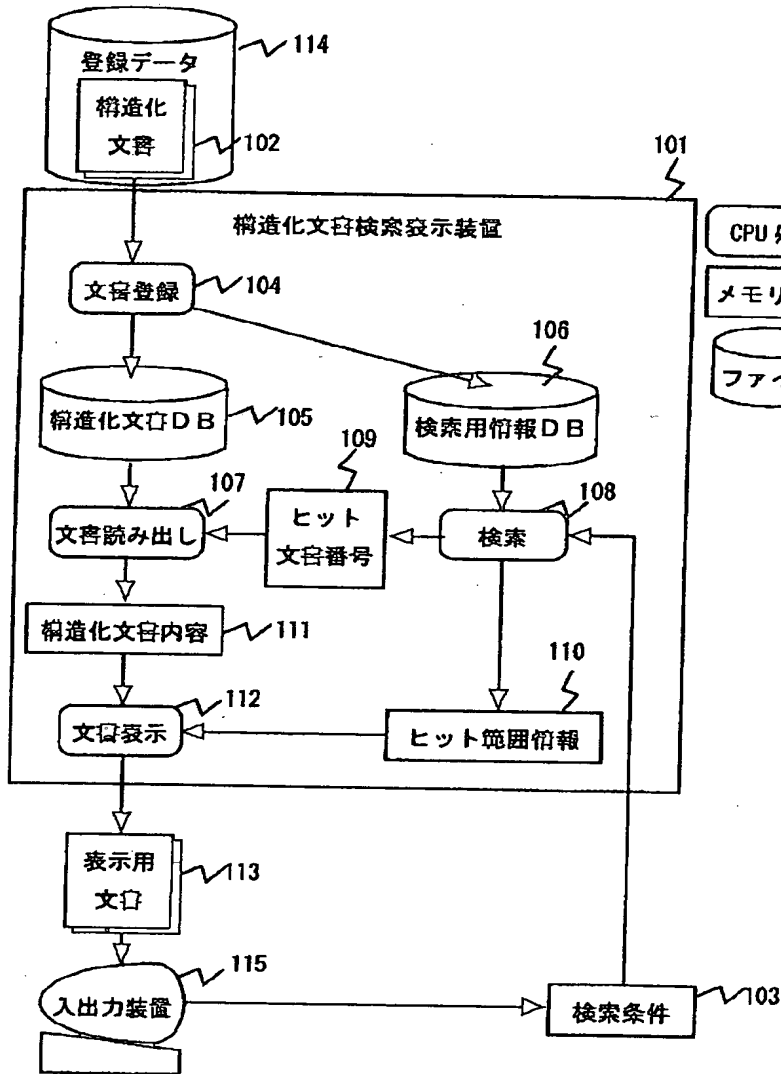
2707 HTML文書

3200 ハイライトタグ文字格納領域の構造
 3400 HTML 文書例
 3401 HTML 文書の表示画面例

3500 ハイライトタグ挿入後HTML 文書例
 3501 ハイライトタグ挿入後表示画面例

【 図1 】

【 図1 】



<実施例1、2の処理ブロック図>

【 図5 】

【 図5 】

<検索用のテキスト>

相違ID	開始文字位置	文字列長さ
2	0	10
5	10	6
6	16	44
:	:	:

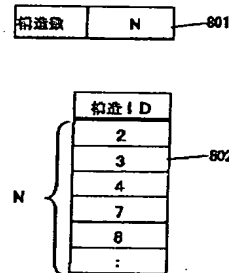
<文字列>

相違化文書相違化文書とは、.....

【 図8 】

【 図8 】

<相違指定の解析結果>



【 図14 】

【 図14 】

<ヒット範囲情報>

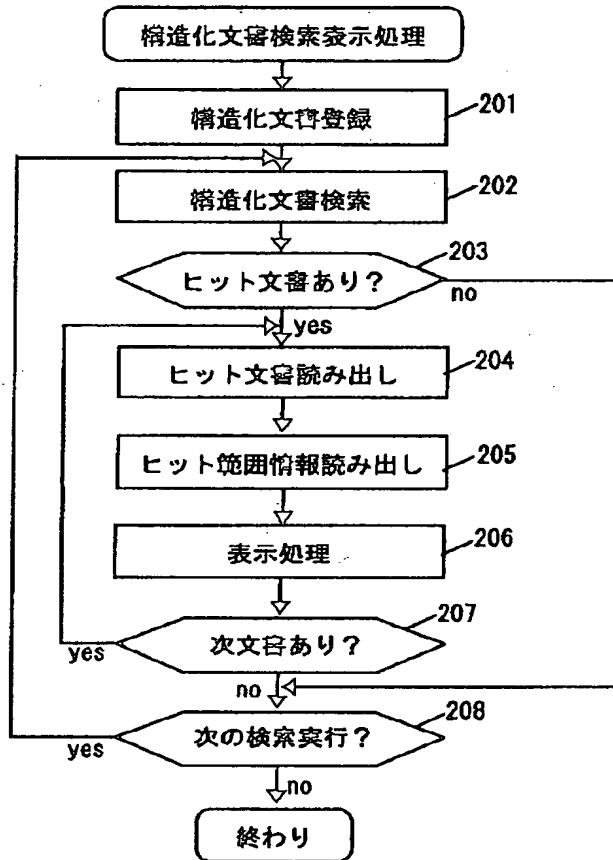
正規化後 相違ID	ヒット範囲	
	開始位置	長さ
2	0	10
9	0	10

正規化後の情報への変換

正規化前 相違ID	ヒット範囲	
	開始位置	長さ
2	0	10
5	0	6
6	0	4

【図2】

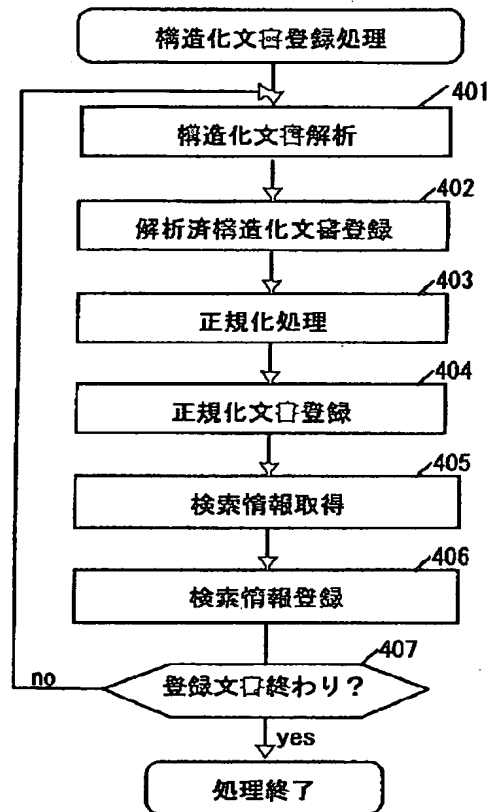
【図2】



<実施例1の文書表示処理のフローチャート>

【図4】

【図4】



<構造化文書登録処理のフローチャート>

【図17】

【図17】

<ヒット範囲情報>

初段ID	ヒット条件	ヒット範囲	
		開始位置	長さ
2	ターム	0	10
9	頻度	0	10
9	初段	0	50
9	ターム	30	4
9	ターム	36	4
9	頻度	30	10

【図18】

【図18】

<ハイライト方法定義>

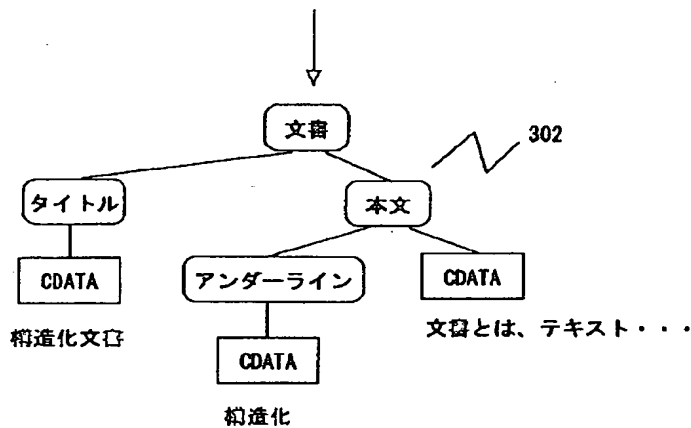
ヒット条件	ハイライト方法	隠蔽初段
ターム	赤色	1
頻度	フォント大	1
距離	アンダーライン	2
初段	例体	3

【 図3 】

【図3】
 <構造化文書の登録>

```

<!DOCTYPE 文書 SYSTEM "文書.dtd">
<文書>
<タイトル>構造化文書</タイトル>
<本文>
<アンダーライン>構造化</アンダーライン>文書とは、テキスト・・・
</本文>
</文書>
  
```



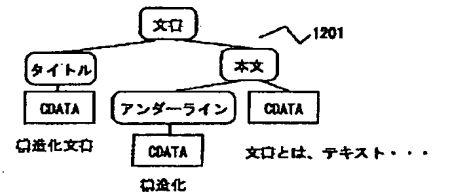
テーブル形式の
データへの変換

構造化識別子	構造化種別	タグ	内容
0	文書	文書	下位構造: 1, 3
1	要素	タイトル	下位構造: 2
2	CDATA	—	文字列: "構造化文書"
3	要素	本文	下位構造: 4, 6
4	要素	アンダーライン	下位要素: 5
5	CDATA	—	文字列: "構造化"
6	CDATA	—	文字列: "文書とは、..."

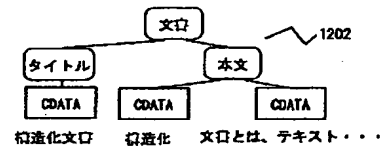
【 図12 】

【図12】

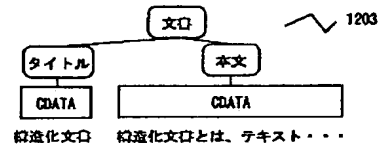
<構造化文書検査用正規化処理>



不要な構造 (アンダーライン) の削除

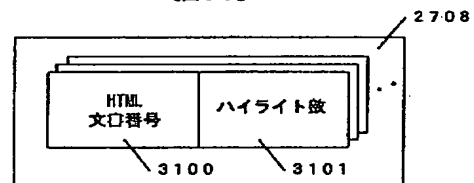


テキストの結合



【 図31 】

【図31】



【 図22 】

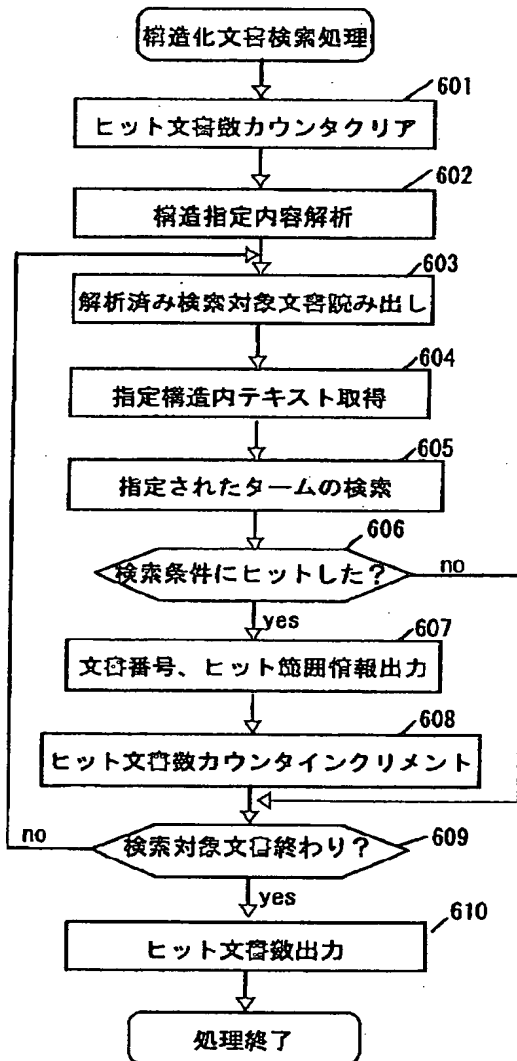
【図22】

<ハイライト表示の例>

構造化文書とは、テキスト中にタグを挿入した文書・・・

【 図6 】

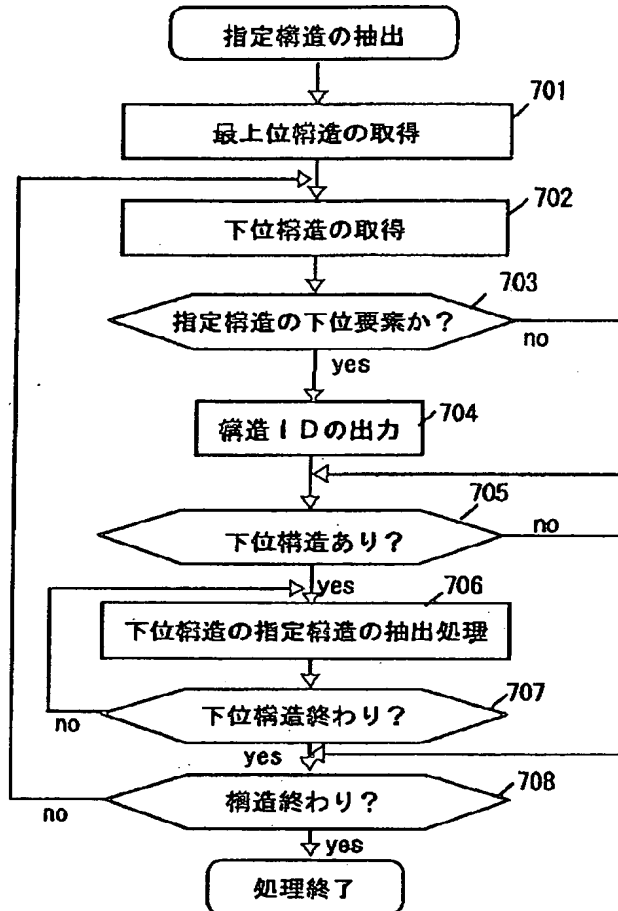
【図6】



<構造化文書検索処理のフローチャート>

【 図7 】

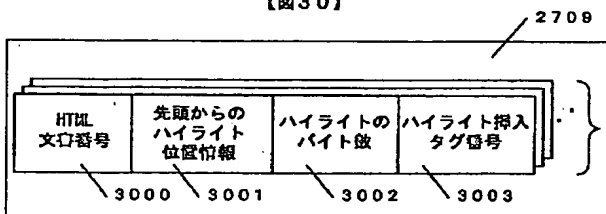
【図7】



<指定構造抽出処理のフローチャート>

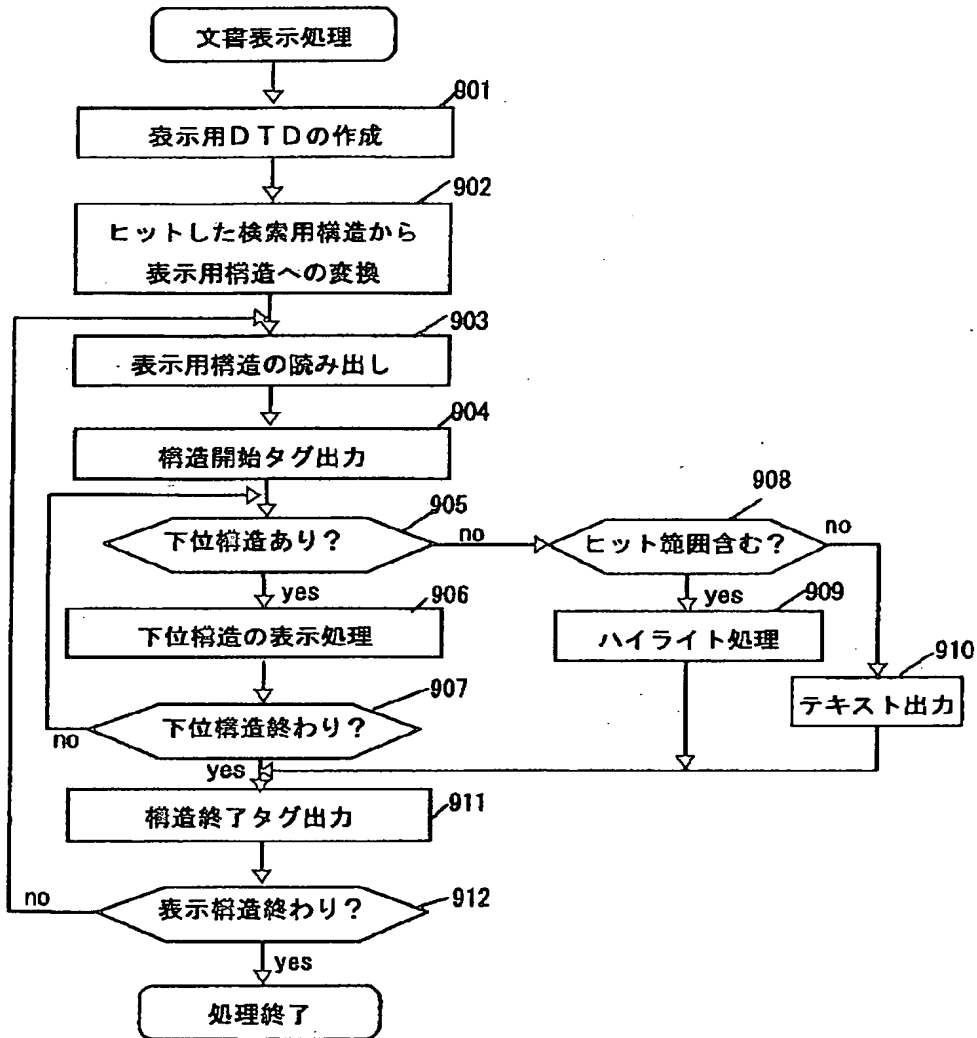
【 図30 】

【図30】



【 図9 】

【 図9 】



<構造化文書表示処理のフローチャート>

【 図43 】

【 図43 】

<文書. タイトル>: "構造化文書" (青色) AND
 <文書. 本文> [並下線]: ("構造化文書" OR C<10 (反転) ("タグ", "挿入"))

【 図10 】

【図10】

<構造化文書の例>

[DTD] (文書. dtd)

1001

```

<!ELEMENT 文書      -- (タイトル, 本文)>
<!ELEMENT タイトル  -- CDATA>
<!ELEMENT 本文      -- (#PCDATA|アンダーライン)*>
<!ELEMENT アンダーライン -- CDATA>

```

[文書インスタンス]

1002

```

<!DOCTYPE 文書 SYSTEM "文書. dtd">
<文書>
<タイトル>構造化文書</タイトル>
<本文>
<アンダーライン>構造化</アンダーライン>文書とは、テキスト、、、
</本文>
</文書>

```



“構造化文書”を検索、結果表示

<ハイライト処理後の構造化文書の例>

[DTD] (表示. dtd)

1003

```

<!ELEMENT 文書      -- (タイトル, 本文)>
<!ELEMENT タイトル  -- (#PCDATA|ハイライト)*>
<!ELEMENT 本文      -- (#PCDATA|アンダーライン|ハイライト)*>
<!ELEMENT アンダーライン -- (#PCDATA|ハイライト)*>
<!ELEMENT ハイライト -- (#PCDATA)>

```

[文書インスタンス]

1004

```


<!DOCTYPE 文書 SYSTEM "表示. dtd">
<文書>
<タイトル><ハイライト>構造化文書</ハイライト></タイトル>
<本文>
<アンダーライン><ハイライト>構造化</ハイライト></アンダーライン>
<ハイライト>文書</ハイライト>とは、テキスト、、、
</本文>
</文書>

```

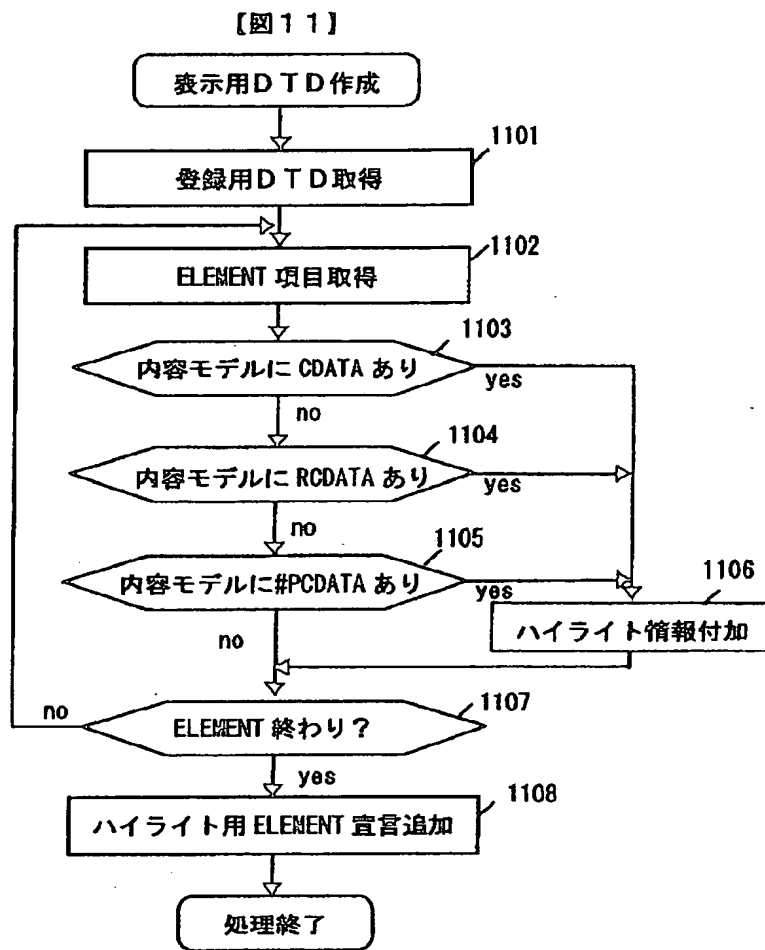
【 図47 】

【図47】

<ハイライト表示の例>

構造化文書とは、テキスト中に  した文書・・・

【 図11 】



<表示用DTD作成のフローチャート>

【 図44 】

【図44】

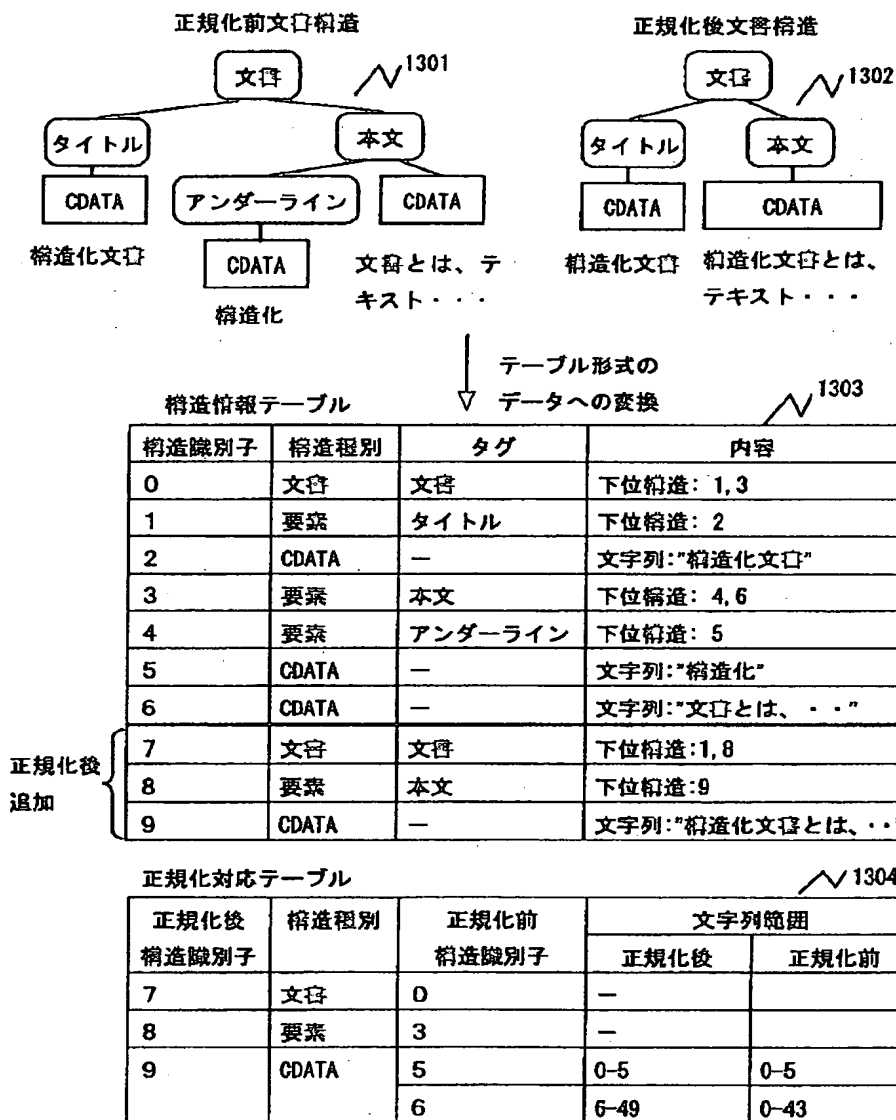
<ヒット範囲情報>

構造ID	ハイライト 方法	ヒット条件	ヒット範囲	
			開始位置	長さ
2	4403 色	ターム	0	10
9	—	属性	0	10
9	4402 下線	構造	0	50
9	—	ターム	30	4
9	—	ターム	36	4
9	4401 反白	属性	30	10

【 図13 】

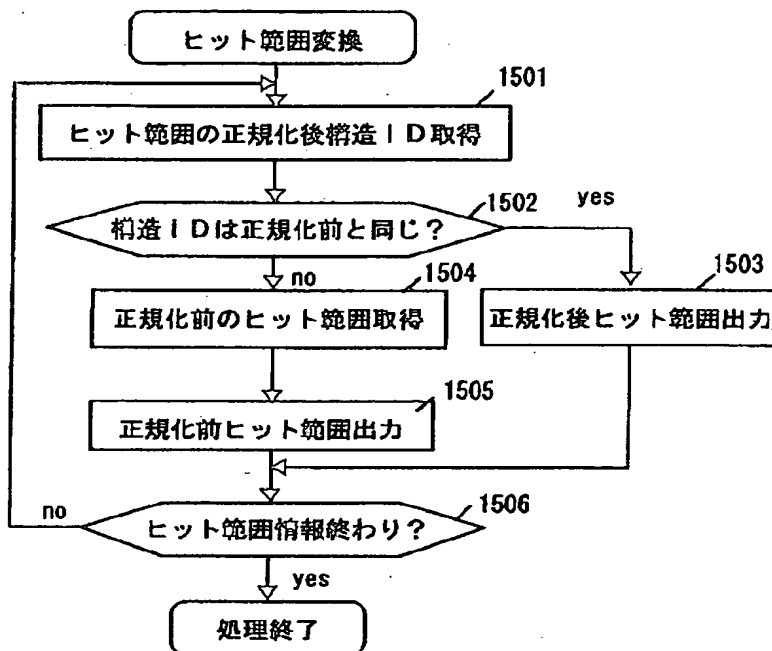
【 図13 】

<正規化文書のテーブル形式のデータへの変換>



【 図15 】

【図15】



<ヒット範囲取得処理のフローチャート>

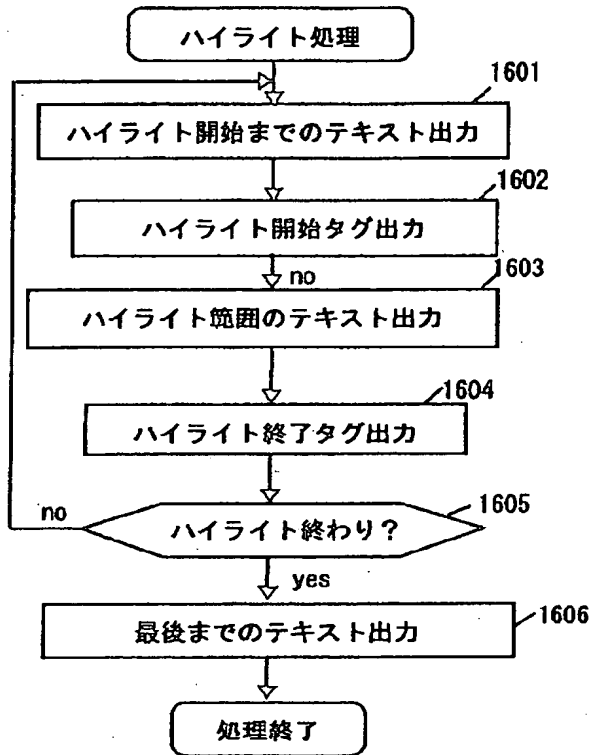
【 図21 】

【図21】

<文符>
 <タイトル><ターム>構造化文符</ターム></タイトル>
 <本文>
 <アンダーライン><構造><頻度>構造化</頻度></構造></アンダーライン><構造><頻度>文符</頻度>
 <頻度>とは、テキスト中に<距離><ターム>タグ</ターム>を<ターム>挿入</ターム></距離>した文符・・・</構造>
 </本文>
 </文符>

【図16】

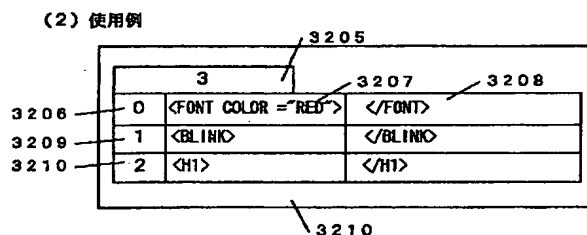
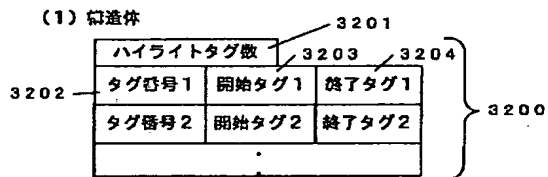
【図16】



<ハイライト処理のフローチャート>

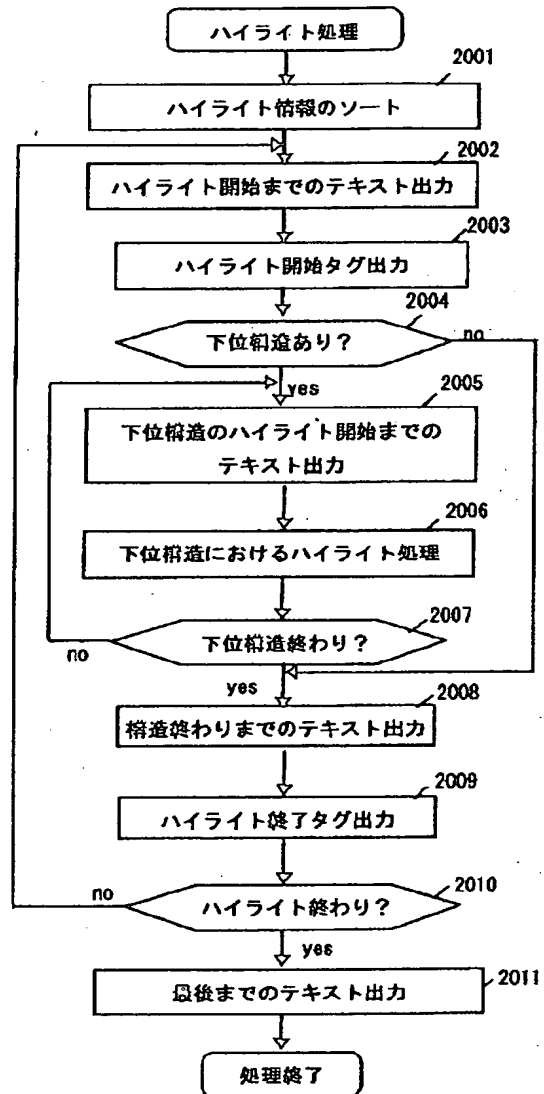
【図32】

【図32】



【図20】

【図20】

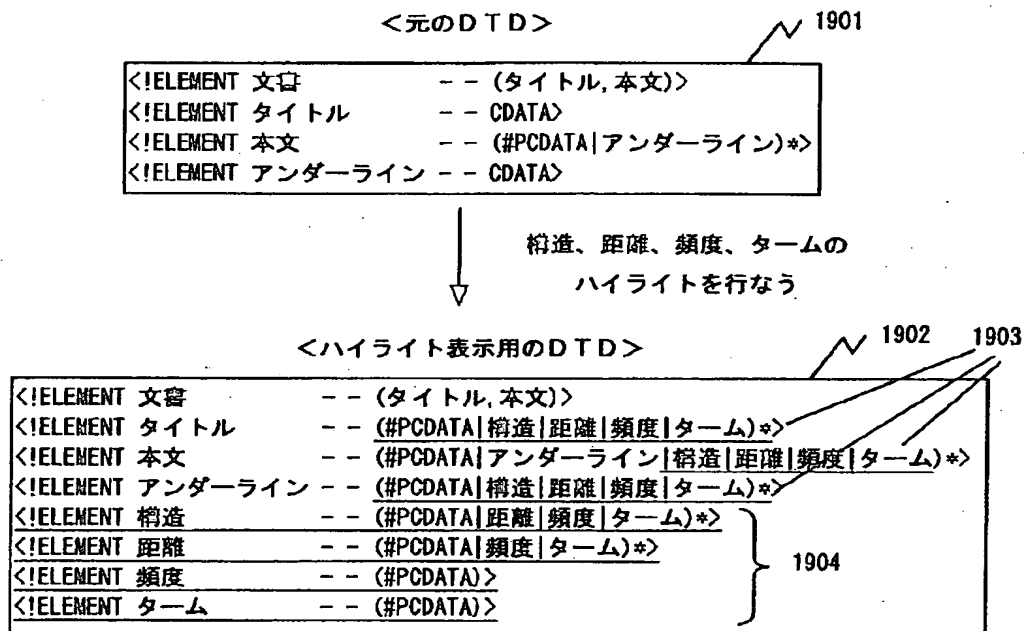


<ヒット範囲取得処理のフローチャート>

【 図19 】

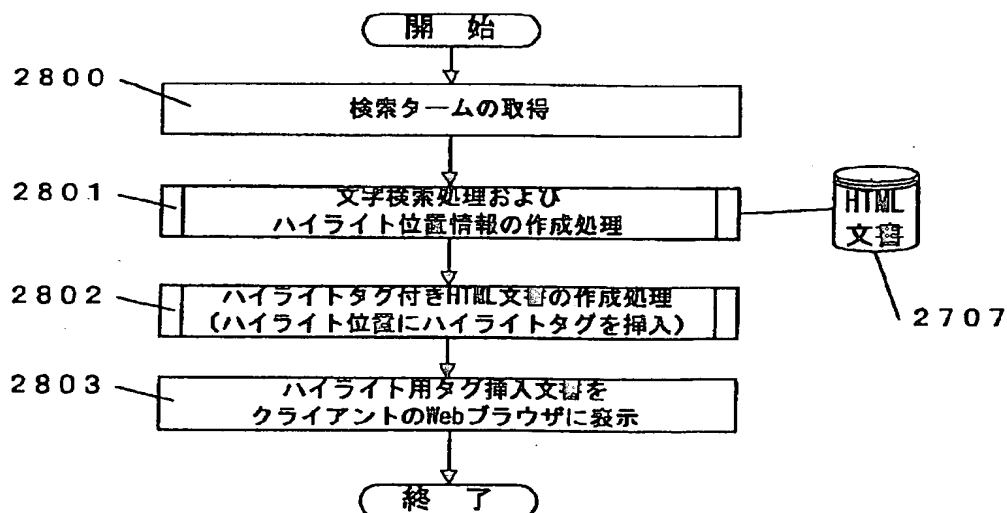
【図19】

<ハイライト表示用DTDへの変換>



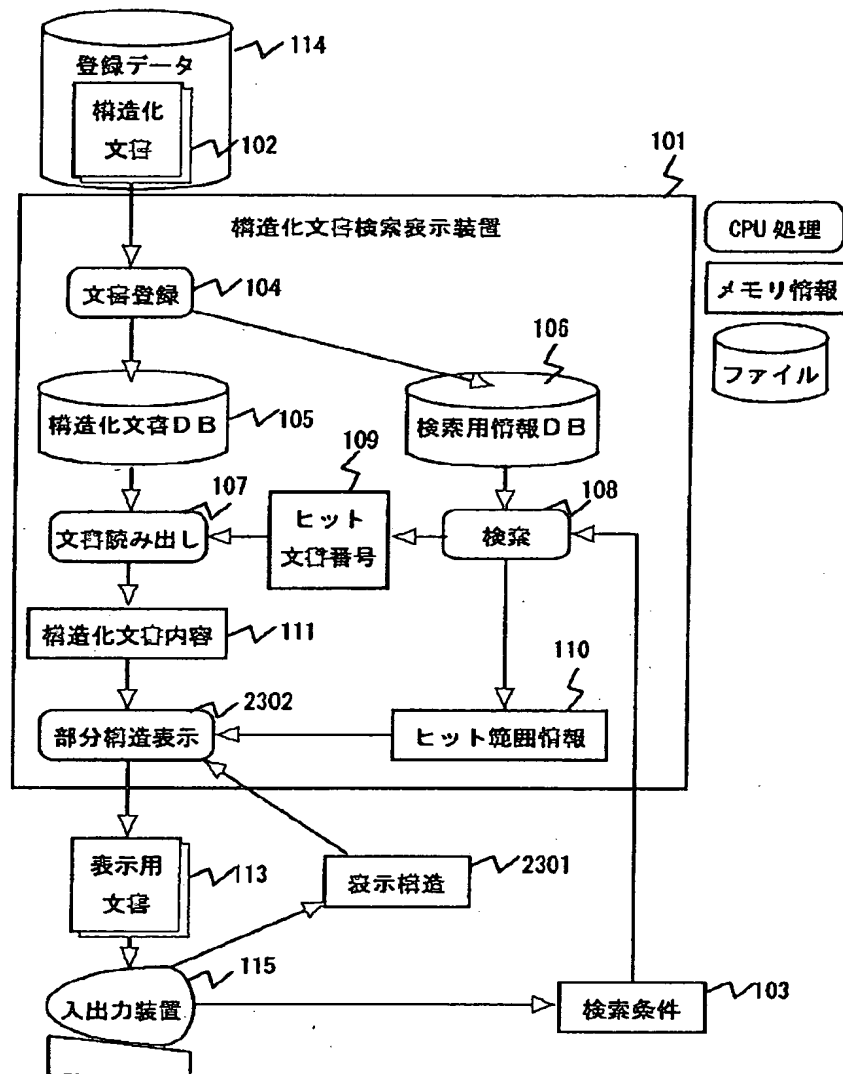
【 図28 】

【図28】



【 図23 】

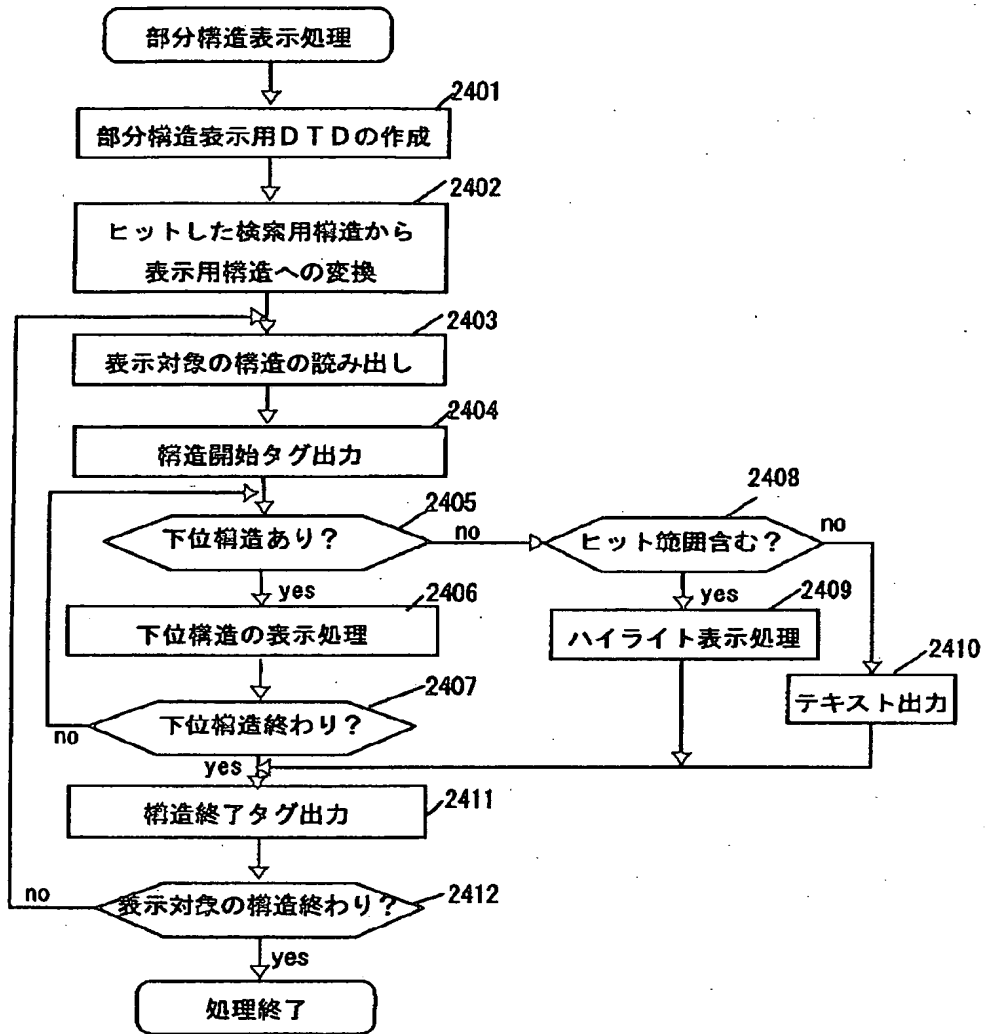
【図23】



<実施例3の処理ブロック図>

【 図24 】

【 図24 】

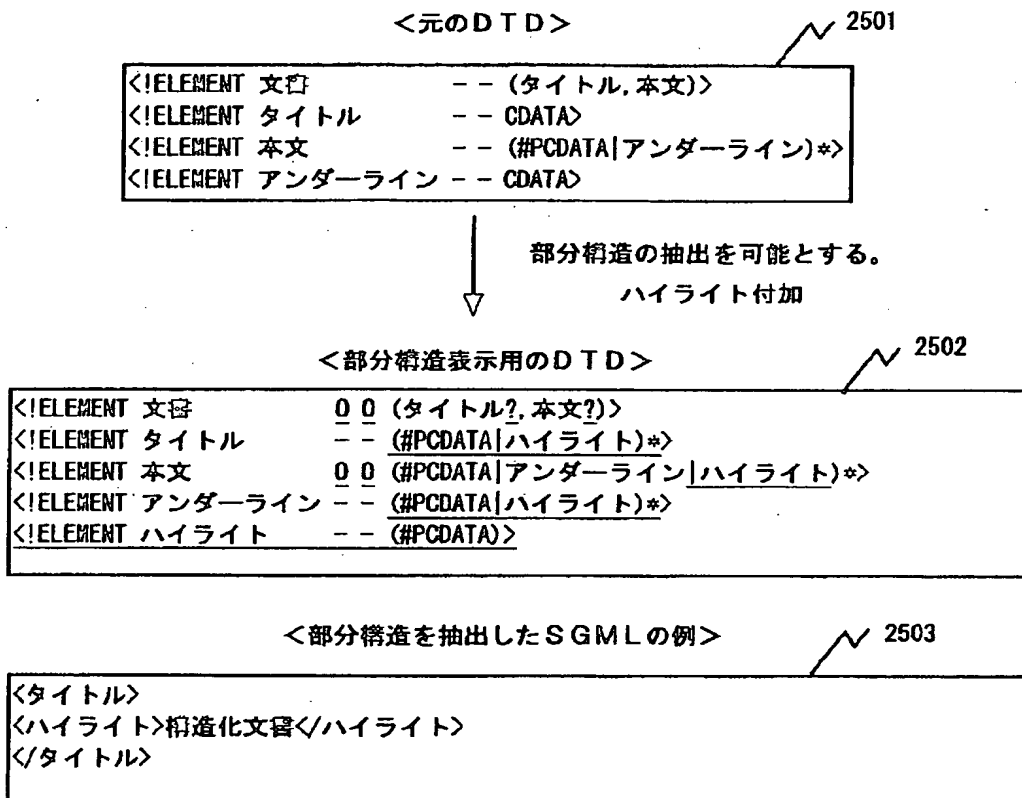


<構造化文書表示処理のフローチャート>

【 図2 5 】

【図2 5】

<部分構造表示用D T Dへの変換>



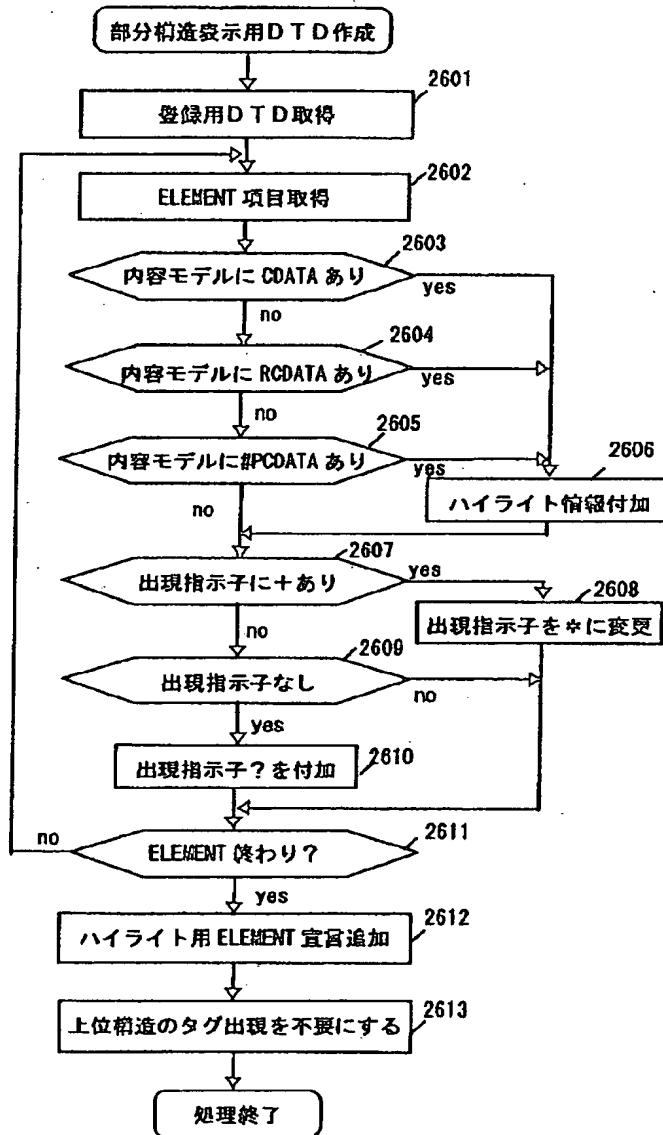
【 図4 6 】

【図4 6】

<文書>
 <タイトル><青色>構造化文書</青色></タイトル>
 <本文>
 <アンダーライン><波下線><フォント大>構造化</フォント大></波下線></アンダーライン><
 波下線><フォント大>文書</フォント大>とは、テキスト中に<反転><赤色>タグ</赤色>を<赤
 色>挿入</赤色></反転>した文書・・・</波下線>
 </本文>
 </文書>

【 図2 6 】

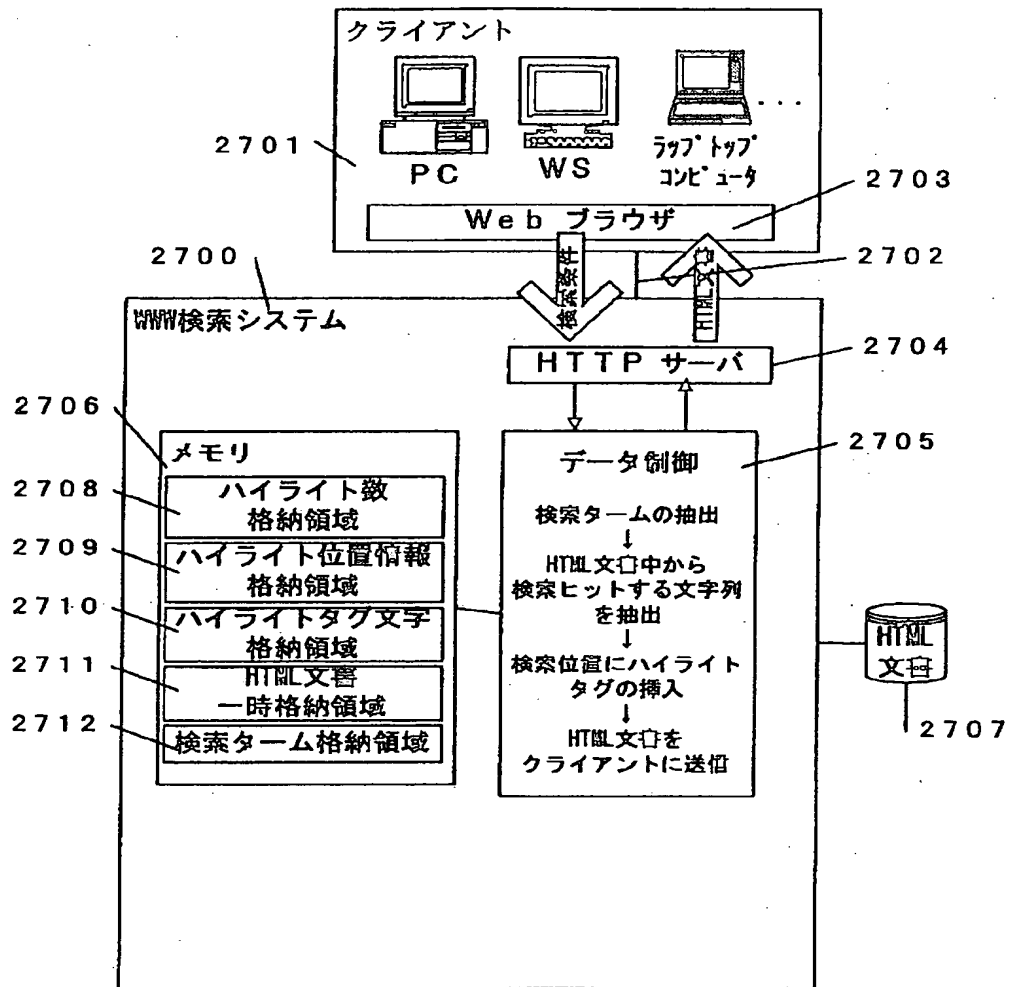
【 図2 6 】



<部分構造表示用DTD作成のフローチャート>

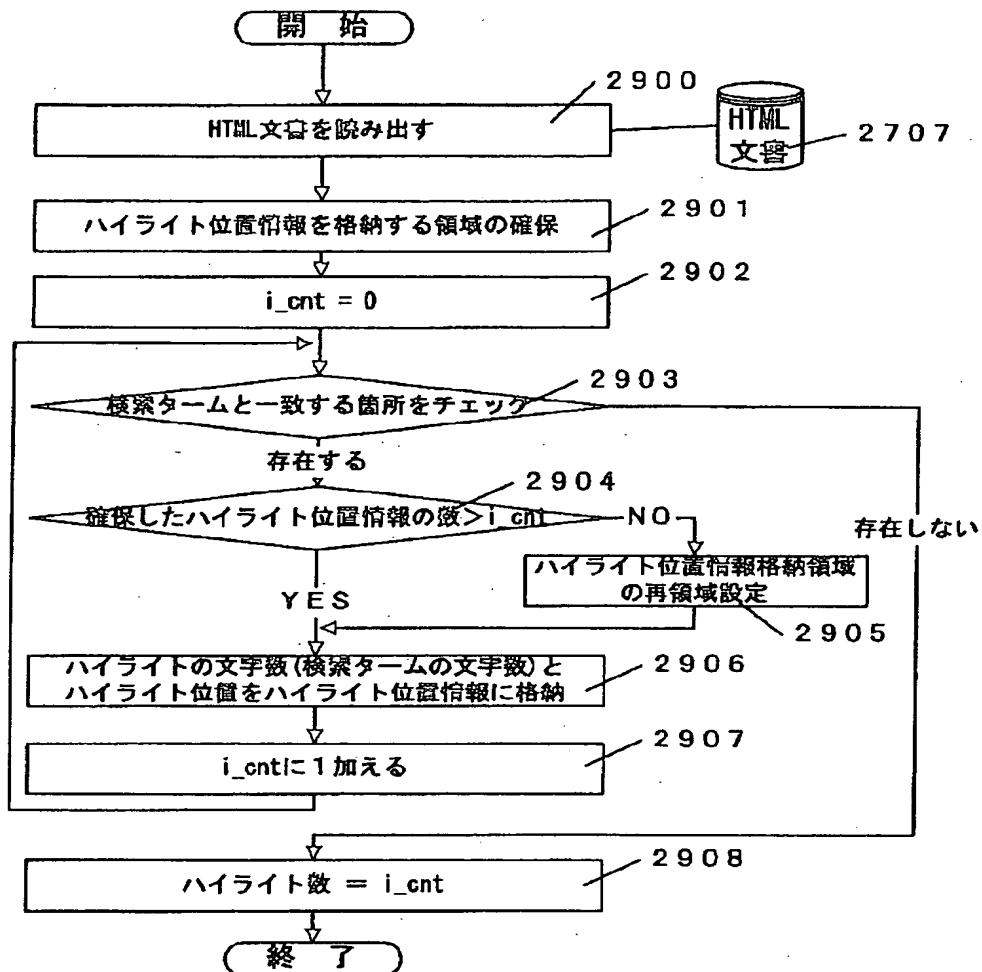
【 図27 】

【 図27 】



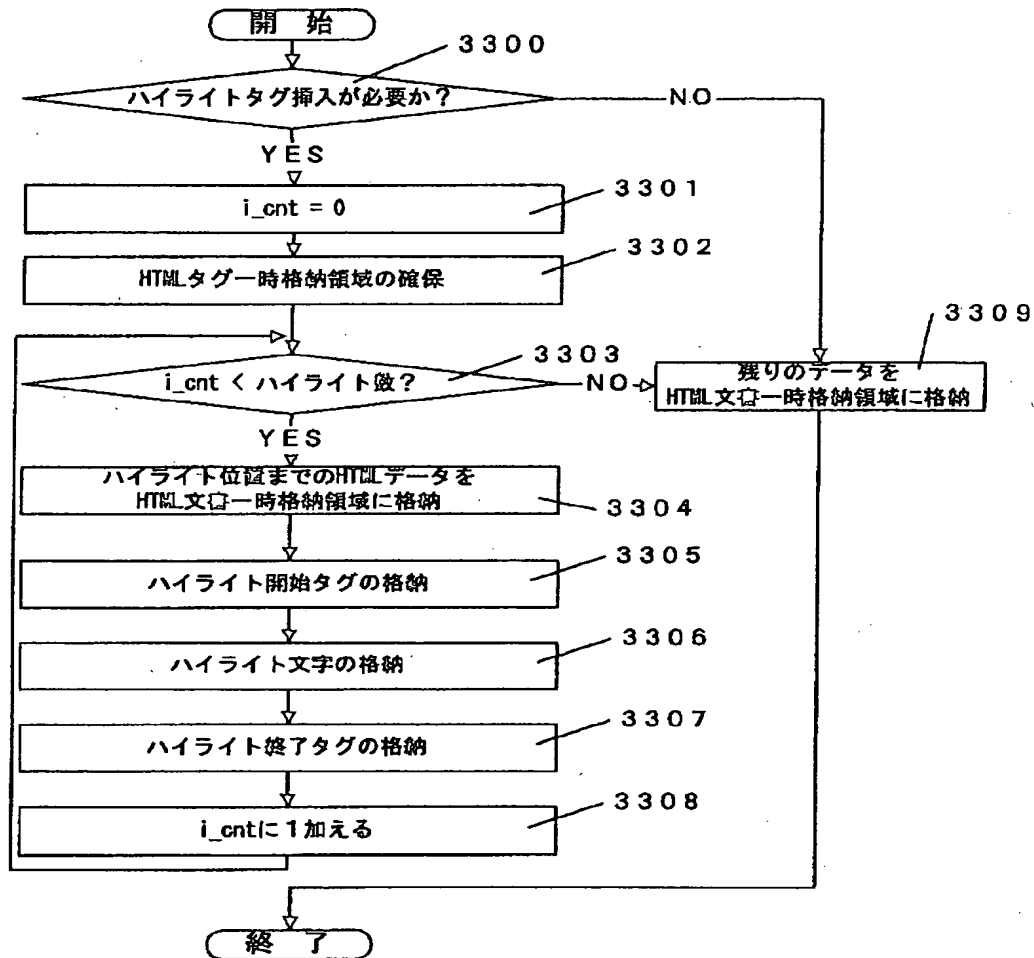
【 図29 】

【 図29 】



【 図33 】

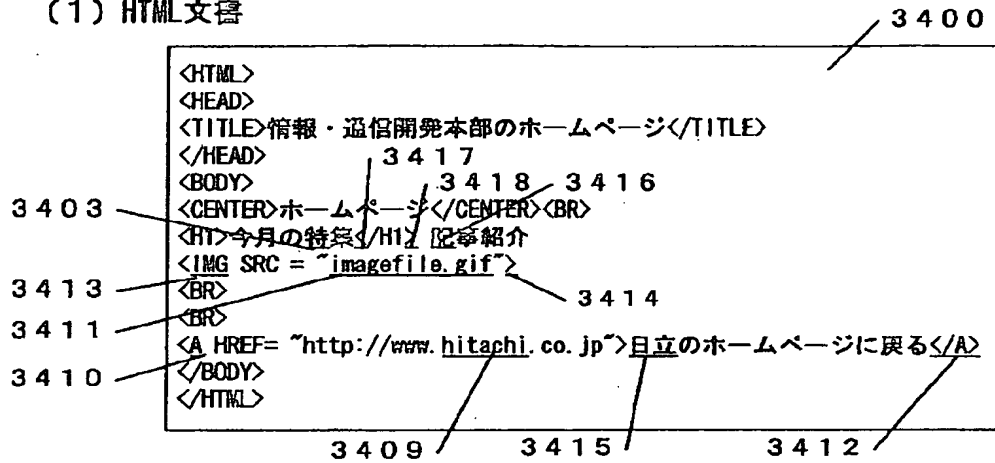
【 図33 】



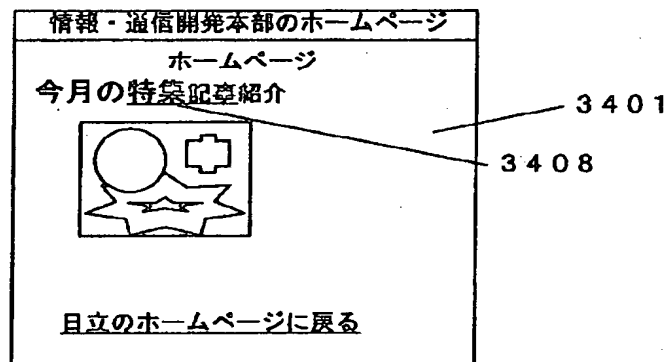
【 図34 】

【 図34 】

(1) HTML文書



(2) 表示画面



(3) ハイライト位置情報格納領域

HTML 文書番号	先頭からの ハイライト 位置情報	ハイライトの バイト数	ハイライト挿入 タグ番号	3402
001	122	4	1	
3404	3405	3406	3407	

【 図35 】

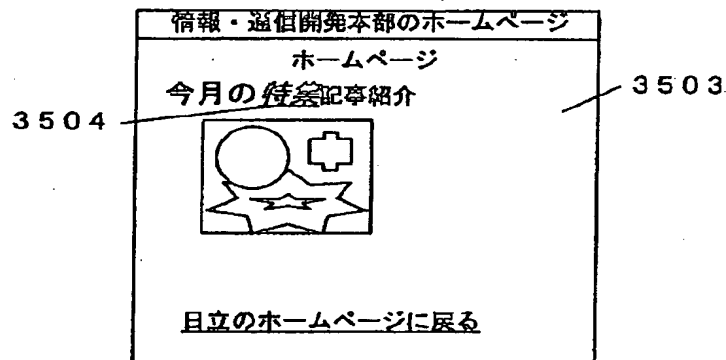
【 図35 】

(1) ハイライトタグ挿入後HTML文書

3500

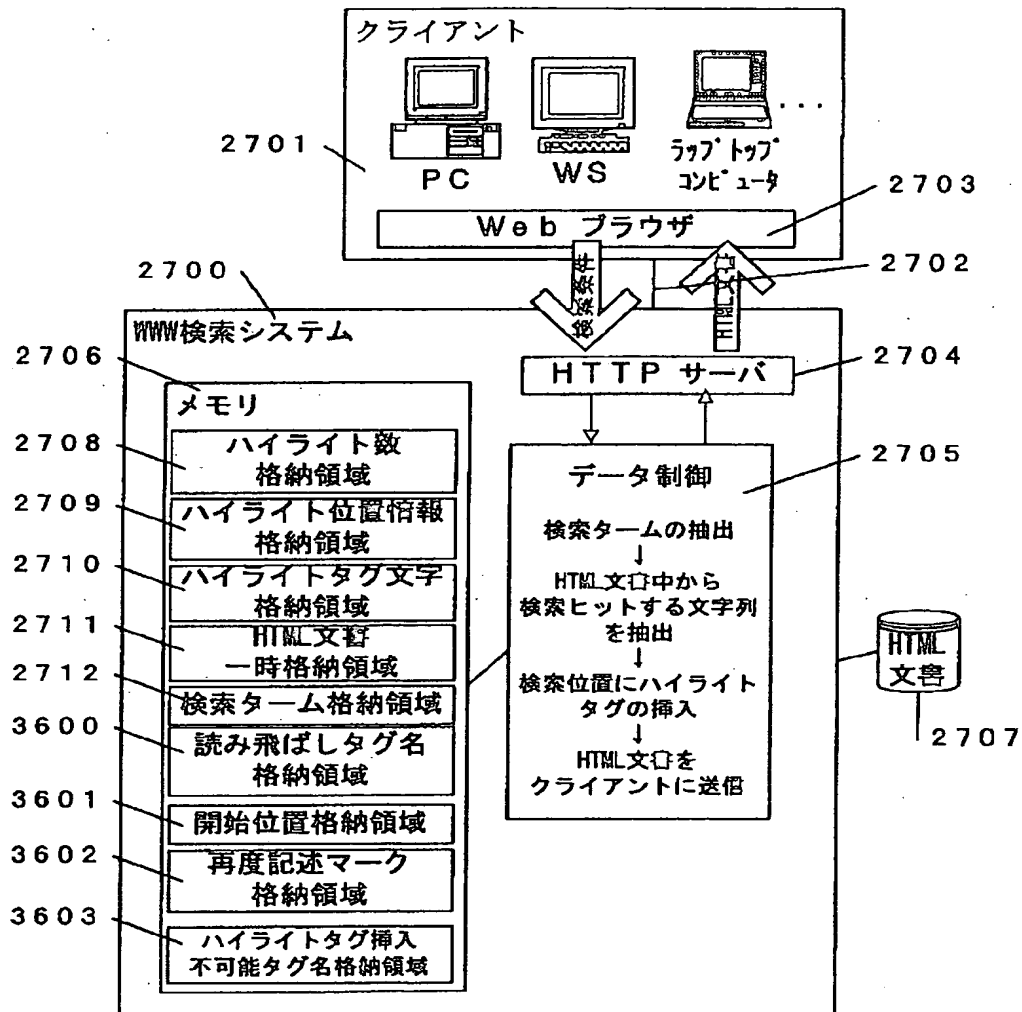
```
<HTML>
<HEAD>
<TITLE>情報・通信開発本部のホームページ</TITLE>
</HEAD>
<BODY>
<CENTER>ホームページ</CENTER><BR>
<H1>今月の<BLINK>特集</BLINK></H1>記事紹介
<IMG SRC = "imagefile.gif">
3501 <BR>
3502 <BR>
<A HREF= "http://www.hitachi.co.jp">日立のホームページに戻る
</A>
</BODY>
</HTML>
```

(2) ハイライトタグ挿入後表示画面



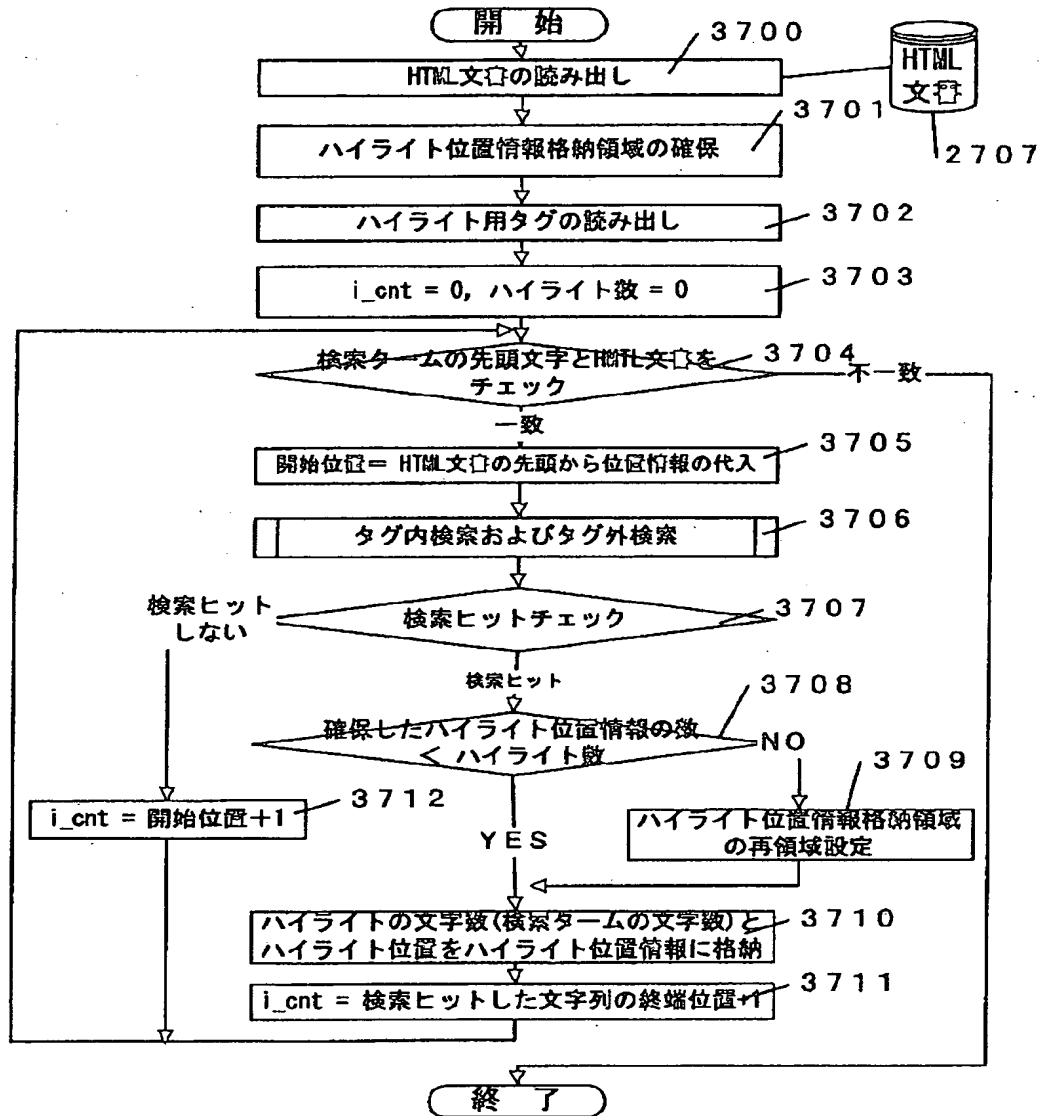
【 図36 】

【 図36 】



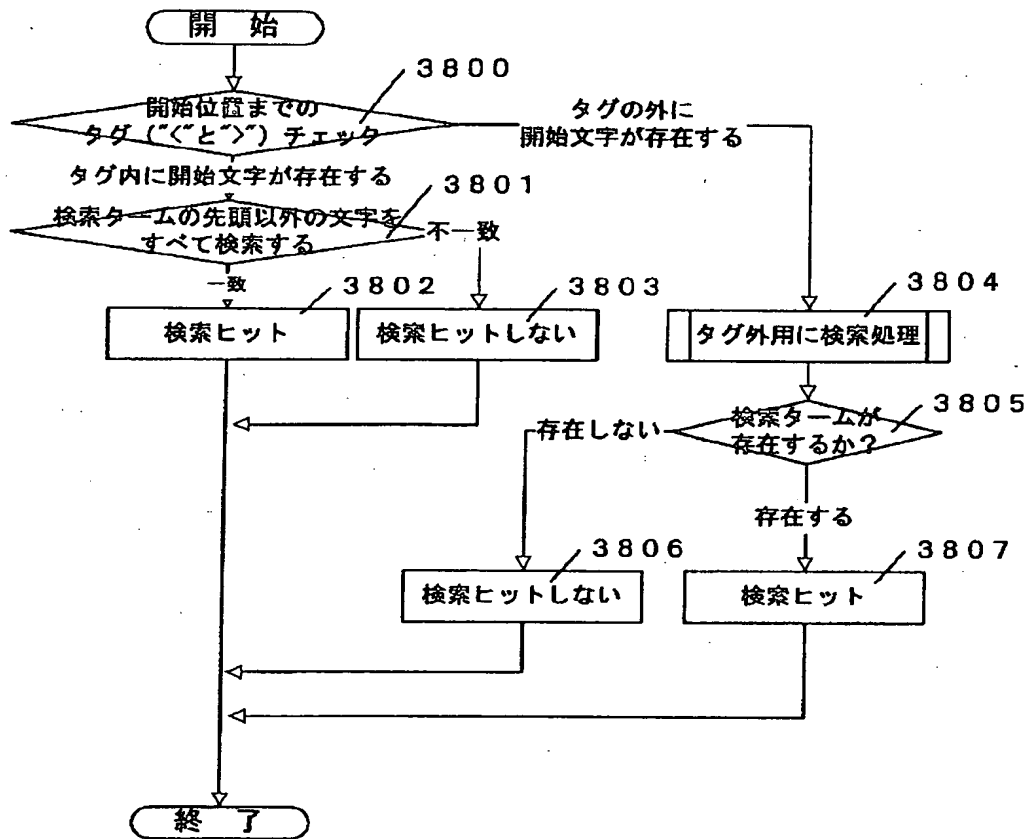
【図37】

【図37】



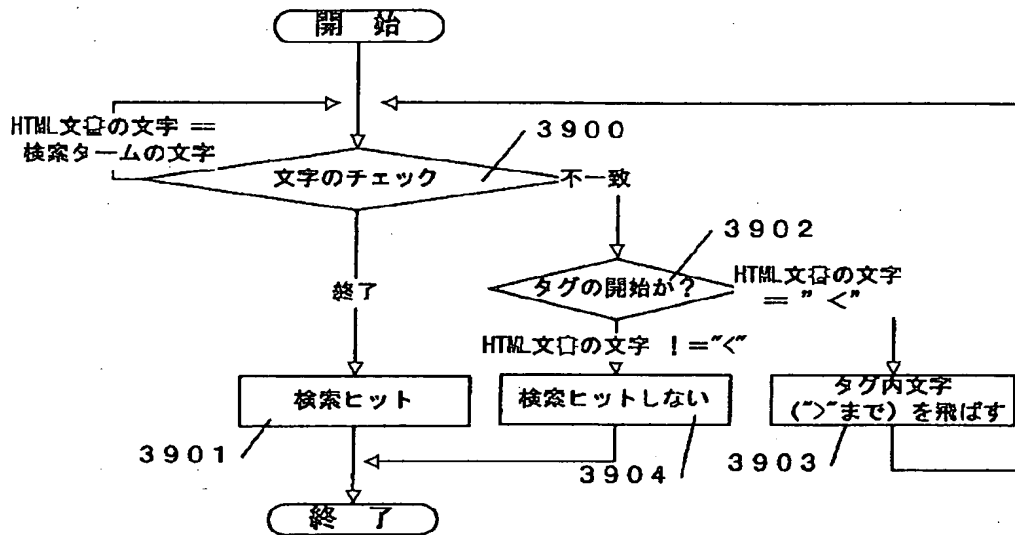
【 図38 】

【 図38 】



【 図39 】

【 図39 】



【 図45 】

【 図45 】

<実施例6におけるハイライト表示用DTDへの変換>

<元のDTD>

```

<!ELEMENT 文書      -- (タイトル, 本文)>
<!ELEMENT タイトル  -- CDATA>
<!ELEMENT 本文      -- (#PCDATA|アンダーライン)*>
<!ELEMENT アンダーライン -- CDATA>
  
```

表示文書用のハイライト情報の付加

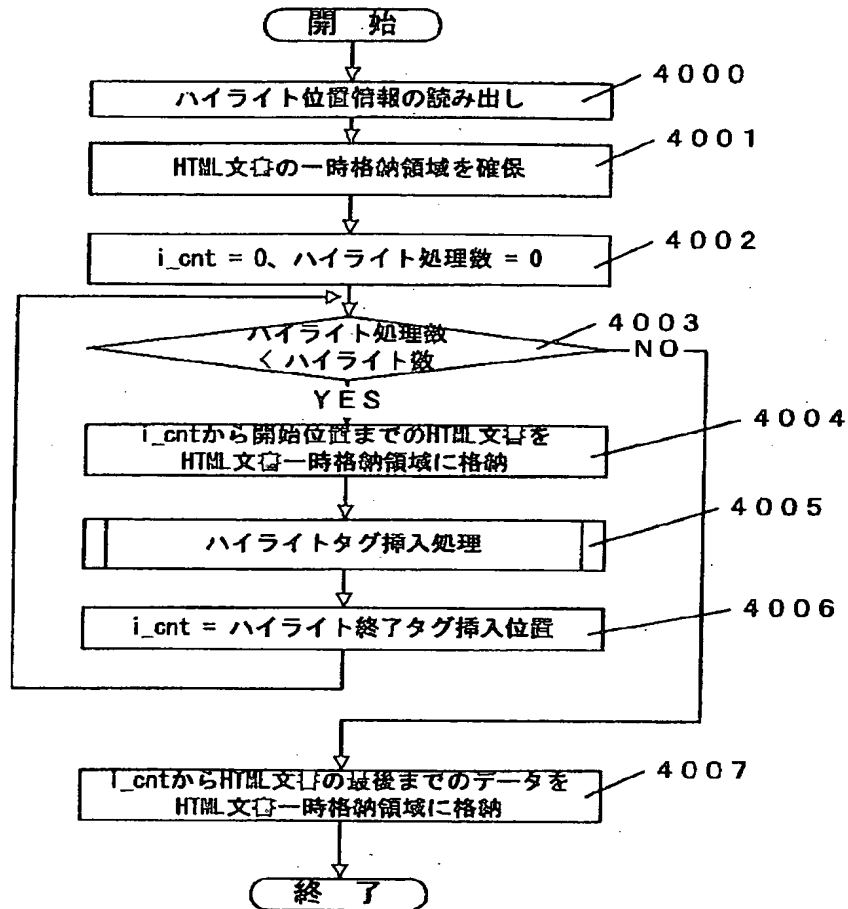
<表示用文書のDTD>

```

<!ELEMENT 文書      -- (タイトル, 本文)>
<!ELEMENT タイトル  -- (#PCDATA|青色)*>
<!ELEMENT 本文      -- (#PCDATA|アンダーライン|波下線|フォント大|反転)*>
<!ELEMENT アンダーライン -- (#PCDATA|フォント大)*>
<!ELEMENT 波下線      -- (#PCDATA|フォント大|反転)*>
<!ELEMENT 反転        -- (#PCDATA|赤色)*>
<!ELEMENT フォント大  -- (#PCDATA)>
<!ELEMENT 青色        -- (#PCDATA)>
<!ELEMENT 赤色        -- (#PCDATA)>
  
```

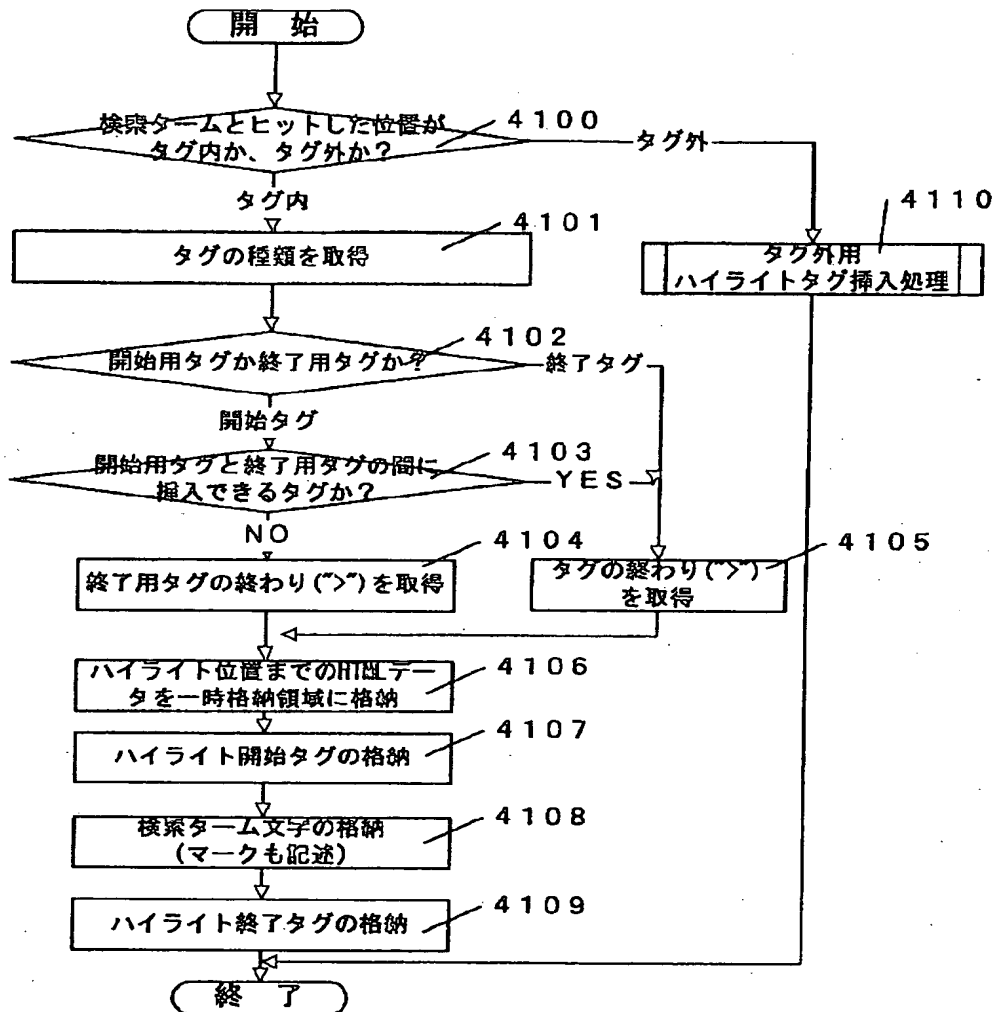
【 図40 】

【 図40 】



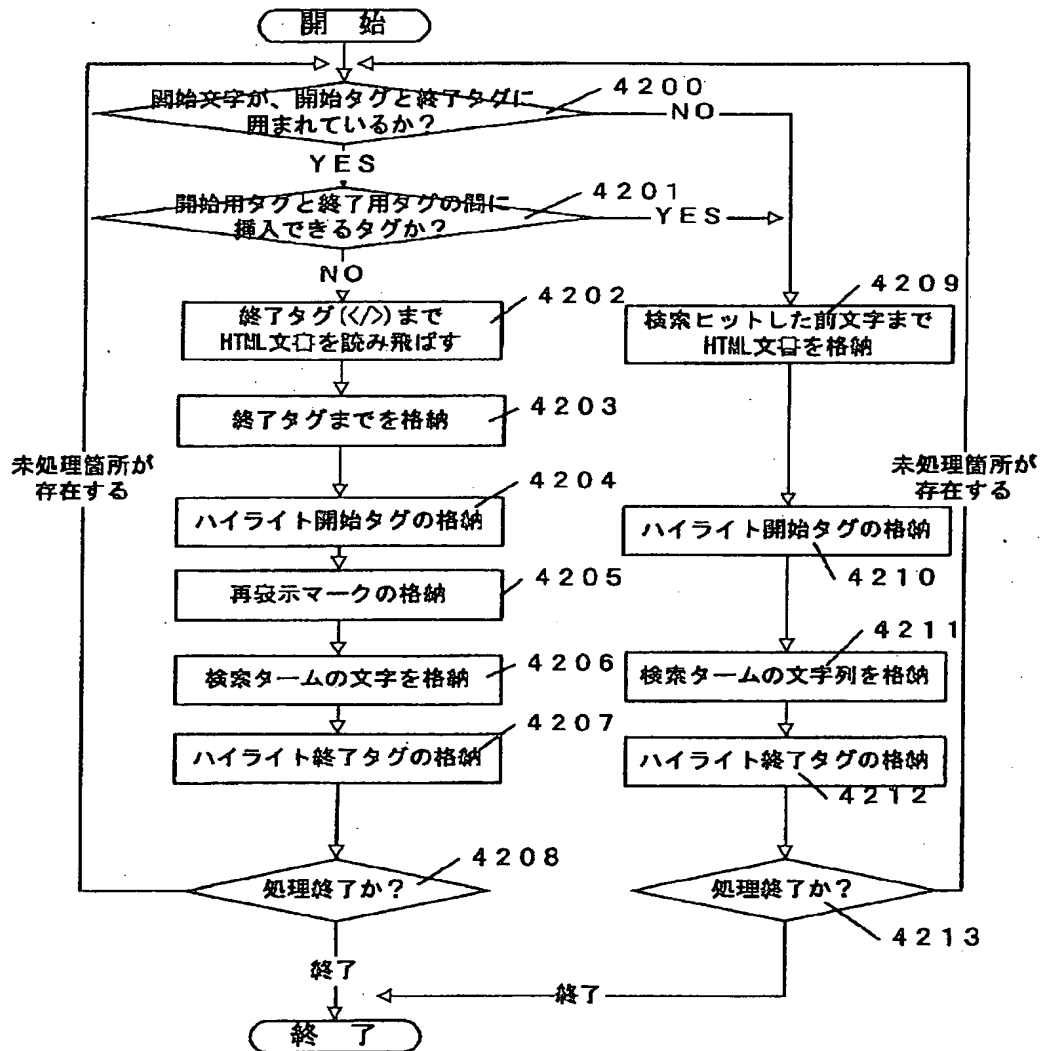
【 図41 】

【 図41 】



【 図42 】

【 図42 】



フロント ページの続き

(51) Int. Cl. 6

識別記号

FI

G06F 15/403

380Z

(72) 発明者 山崎 紀之
 神奈川県戸塚区戸塚町5030番地 株式会社
 日立製作所ソフトウェア開発本部内

(72) 発明者 青山 ゆき
 神奈川県横浜市都筑区加賀原二丁目2番
 株式会社日立製作所システム開発本部内